

MARKOV CHAIN BASED QUERY CLASSIFICATION FOR USER INTENT IMAGE SEARCH ENGINES

THI THI ZIN¹, PYKE TIN¹ AND HIROMITSU HAMA²

¹Faculty of Engineering
University of Miyazaki
1-1, Gakuen Kibanadai-nishi, Miyazaki 889-2192, Japan
thithi@cc.miyazaki-u.ac.jp

²Research Center for Industry Innovation
Osaka City University
3-3-138, Sugimoto Sumiyoshi-ku, Osaka 558-8585, Japan

Received February 2016; accepted May 2016

ABSTRACT. *Today Web image search engines have well recognized that a challenging and difficult problem of determining the user intent of Web searches should be thoroughly investigated. In this aspect we propose a Markov chain query classification method to understand user behavior and intent based on the queries which users used for searches. Specifically, we analyze some samples of queries from different Web search engines and classify the user queries into three classes such as informational, navigational and transactional classes. By assuming these three classes represent the type of contents that a user desired to express his or her intents, we implement the Markov chain based classification process. Some experimental results are shown to capture user intents by using a set of queries submitted to the Web search engines.*

Keywords: Image search engine, Query classification, Markov chain, User intents, Informational, navigational and transactional classes

1. Introduction. Nowadays, a tremendous amount of images are available on the World Wide Web. One of the most desirable things that the users want is to have an effective Web image search engine that returns images of users' needs [1]. In this aspect, we well recognized that understanding user behavior while users are searching images in web-scale image search has played a key role in retrieving user satisfactory images. On the other hand, for establishing a good image search engine the content providers are interested in determining what types of images users want. Moreover, the information search engines have become indispensable tools in our daily life for finding and accessing important information which people need. Taking this fact into account it is worthwhile to look into the ways how people are using the search engines and what types of queries they are mostly using so that we can improve our search engines to provide the information that users required or user intent which users address behind the queries they input to the search engines.

In order to develop a more effective image search engine, the characteristics of users' image queries and text queries have been investigated. Several studies have shown that text query alone is not sufficient to find images what users need [2-6]. It is, therefore, important to analyze image queries from image search engines to understand how image requests differ from the above observations, and to examine the limitations or difficulties faced by Web image search engines in dealing with a variety of image needs.

Most of the research has focused on specific collections or specific groups of users [7,8]. In this paper, we classify user query to understand user intent. In particular a Markov chain classification approach is proposed by using both textual and visual features. To be specific we organize the rest of paper as follows. In Section 2, some related works are

reviewed and in Section 3 the overview of proposed method and some technical details are presented. Some experimental works are shown in Section 4. Finally we conclude the paper in Section 5.

2. Some Related Works. A sizable number of researchers have attempted to discover the intent of Web users by analyzing the types of images they look for and the queries they used while searching [9,10]. In doing so, their works have focused on the classification process by comparing only informational and navigational in order to simplify the problem. They also used supervised and unsupervised learning to classify Web queries. It is well recognized that a comprehensive evaluation of a substantial set of Web searching queries will significantly enhance understanding user intent in information searching [11]. In addition, several research works have been done in the literature for re-ranking the search results returned by existing text-based image search engines and deliver the images which are more relevant to user queries [12,13]. Also a hierarchical clustering technique using visual, textual and link information to derive users' intentions through the queries is proposed in [14,15].

According to existing literature, efforts at classification of Web queries have usually involved small quantities of queries manually classified [16-18]. There has been little effort on automated classification of queries for user intent. It is these issues that motivate our research. A comprehensive evaluation of a substantial set of Web searching queries will significantly enhance understanding user intent in Web image searching. Therefore, the present considerations of user query classification are very significant to understand user behaviors for Web search engines. Moreover, the use of stochastic concepts such as Markov chain and random walk are strong technical points for the researchers in various research fields. Also, we can understand how the concepts evolved academically since they were first developed by the Russian Mathematician A. A. Markov.

3. Proposed Method for User Query Intent Analysis. The overview of proposed Markov chain based user query intent analysis for image search engine is described in Figure 1. First a seed query is identified from a particular Web search engine such as Google image search, Yahoo image search, and Bing image search. Then for this query the search engine returns thousands of images ranked by the keywords extracted from the surrounding text.

Among them, we first select a few queries, which should be representative queries of the specific domain, and then we map these queries to concepts of the search engine we have used. By using these concept mappings, we can classify the categories to which the queries belong, their neighboring concepts and the concepts they link to in the search

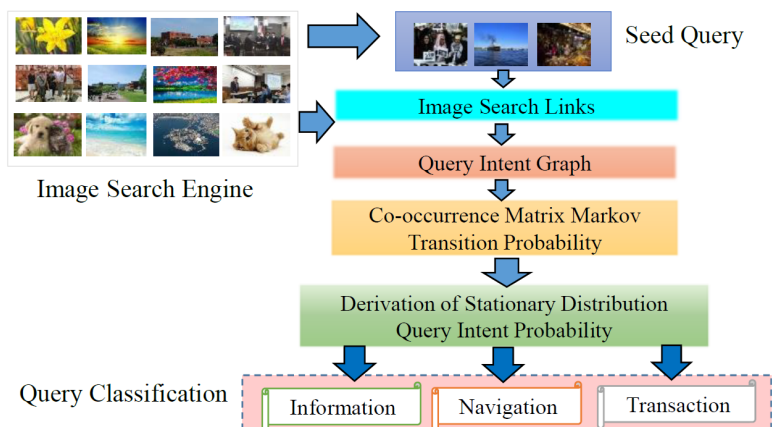


FIGURE 1. Overview of proposed query intent classification system

engine. Thus, we can easily collect a large amount of representatives for search engine concepts and categories as a set of seed concepts.

According to particular search engine concepts, which link to each other through image contents or category links, often share similar topics. To complement the knowledge that is not covered by the seed concepts we collect, we first build a network link graph in which each entry denotes a link relation either within concept image and categories or between image and categories. After that, through a random walk on the built search engine link graph based on the intent label vector, we can get a vector of probabilities measuring how probable each concept in the link graph belongs to the defined intent.

For a user's input query, if it is already covered by search engine concepts, we can easily judge whether the query has the defined intent or not based on the intent probability of the matched search engine concepts. For a query that is not covered by search engine concepts, we employ a method similar to explicit semantic analysis to map the query to its most related search engine concepts and make the intent judgment based on the intent probability of the mapped search engine concepts.

3.1. Search engine link graph construction. Based on the query image and category links of the search engine, we can construct a link graph $G = (X, E)$ where $X = I \cup C$ where I represents a sequence of concept images $I_1, I_2, \dots, I_{m-1}, I_m$ and C represents a set of search engines categories $C_1, C_2, C_3, \dots, C_n$. Each edge in E connects a node within concept images, within categories, or between concept images and categories. There are edges between two nodes in the same set and also between nodes in the different sets. We also note that there are some redundant links between images, and the existence of a link does not always imply that the two images are related. This is due to the fact that many surrounding texts in an image link to other images just because there are entries for the corresponding texts. Therefore, to assure topic relatedness of linked entries, two entries have an edge only when they link to each other in the search engine.

Let W represent an $(m+n) \times (m+n)$ weight matrix, in which element w_{ij} equals the link count associating nodes between x_i and x_j in the matrix. Since the link between two nodes is undirected, the matrix W is symmetric due to $w_{ij} = w_{ji}$. Furthermore, we assume that there is a small set of seed concepts or categories, denoted as X_L , that is manually selected from search engine as positive examples with respect to a specific intent. Given the constructed link graph and the labeled set $X_L = \{X_{L1}, \dots, X_{Lp}\}$, our goal is to automatically propagate labels to more concepts and categories in the set X/X_L .

3.2. Random walk on the link graph. Before introducing the algorithm, we first describe the intuition behind it using Markov chain random walk. Consider the example presented in Figure 2. Suppose we labeled four concepts "Research", "Conference", "Research Category", and "Research Paper" as seeds for the research intent. Here we assume that their immediate neighbors also have the same kind of intent to some extent.

Iteratively, we propagate the intent of research from its seed concepts to their neighborhood nodes. In this example, "Presentation", "Conference Advisory", "Research Theme", "Conference Committee" and "Host Country" have a high travel intent probability after propagation. We define transition probabilities $p_{ij}(t)$ as the probability from node i at time t to node j at time $(t+1)$. We then have

$$p_{ij}(t) = w_{ij}(t) / \sum w_{ij}(t) \quad (1)$$

The Markov chain matrix P is given as

$$P = [p_{ij}(t)] \quad (2)$$

From the Markov chain point of view, the n th power of P will give the probability of the transitions from node i to node j in n transitions or n steps. This also gives a measure of the volume of the paths from one node to another. If there are many paths, the transition

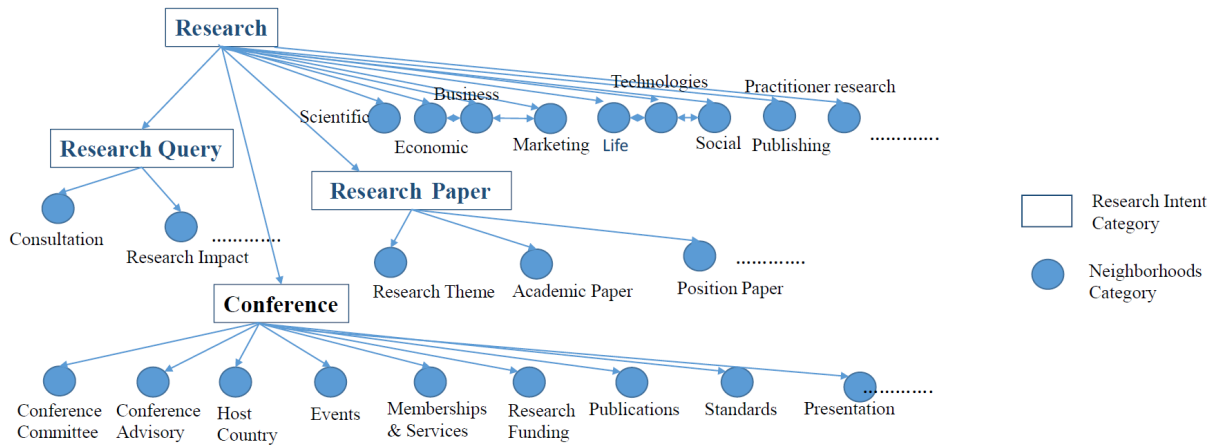


FIGURE 2. Random walk link graph

probability will be higher. Since the matrix P is a stochastic matrix, the largest eigenvalue of P is 1 and all other eigenvalues are in $[0, 1]$. Consequently, it is easy to see that the value in the entry of the stationary distribution π of transition matrix P is the posterior probability ($j = 1, \dots, m + n$) that the node x_j is associated with a specified intent. Therefore, each search engine concept or category is assigned a probability, i.e., $\pi(x_j)$, reflecting the degree of intent.

3.3. User intent classification process. We define here an initial taxonomy of user intent in an image search engine as:

1) *Informational*: That is defined as a user searching for an image to obtain new knowledge or to achieve the answer to a question by looking at an image and by extracting new knowledge from it;

2) *Navigational*: This is a special instance where a user is navigating to find a photo or an image she knows it exists, but does not know how the content of the photo or image exactly looks like. Another aspect can be a class of images in which a user is searching for a specific photo, which he has already in mind and knows how the content of the photo looks like;

3) *Transactional*: This is a case where a user is searching for an image which she wants to obtain for further use.

In this paper we employ the classification process by using the extended taxonomy features. We assume that there can be the overlap of initial classes such as an overlap of the informational and navigational classes as well as between the navigational and transactional classes. For example, the first navigational intent task – “Imagine a friend of you bought a new car which you have not seen yet. He took a photo of that car and uploaded it to Flickr. Find out how your friend’s car looks like.” – can also be classified as an informational class task, since also knowledge is obtained when a user looks at the photos of the car that he has not seen yet. Due to the fact that a task may be assigned to more than one class, the expression “taxonomy” is not valid anymore: the initial taxonomy would assign an intention – in the context of user intent classification – to one single class. This is no longer true for our classification scheme.

Another factor needs to be considered is user search time. There are two different time axes which have to be taken into account: one axis describes one single search session of users and the development of their intent over this search session; the second time axis could describe several search sessions of the same users and the development of their intent over several search sessions.

TABLE 1. Co-occurrence matrix for user intent classes

	Q1	Q2	Q3	Q4	Q5	Q6	Q7	
Q1	52	62	47	46	32	34	25	Q1 Knowledge Intent
Q2	62	84	80	51	36	49	41	Q2 Navigation Intent
Q3	47	80	90	45	41	42	46	Q3 Navigation and Knowledge Intent
Q4	46	51	45	40	40	35	32	Q4 Imagination Image Intent
Q5	32	36	41	40	12	20	10	Q5 Transaction: Potential Future Use Intent
Q6	34	49	42	35	20	28	15	Q6 Knowledge and Imaginary Image Intent
Q7	25	41	46	32	10	15	16	Q7 Transaction and Imaginary Image Intent

TABLE 2. Markov chain transition matrix

0.174497	0.208054	0.157718	0.154362	0.107383	0.114094	0.083893
0.153846	0.208437	0.198511	0.126551	0.08933	0.121588	0.101737
0.120205	0.204604	0.230179	0.11509	0.104859	0.107417	0.117647
0.15917	0.176471	0.155709	0.138408	0.138408	0.121107	0.110727
0.167539	0.188482	0.21466	0.209424	0.062827	0.104712	0.052356
0.152466	0.219731	0.188341	0.156951	0.089686	0.125561	0.067265
0.135135	0.221622	0.248649	0.172973	0.054054	0.081081	0.086486

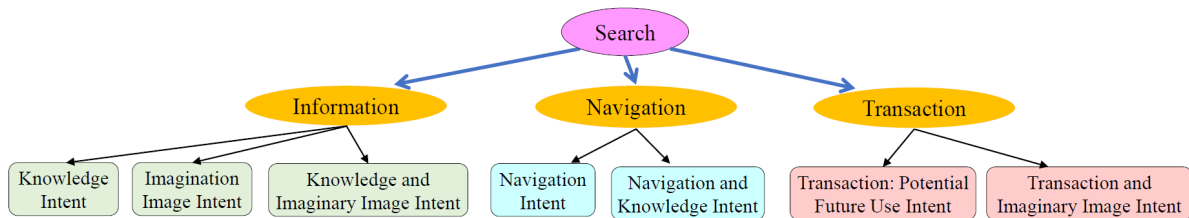


FIGURE 3. Example of query intent classes

4. Simulation Results for Query Classification. In order to illustrate the user query intent classification process proposed in the previous section we generate some image data from Flickr.com. We created a classification process and results by using focus the Flickr Web site. We are able to classify the user’s goal during the search process. Based on the data generated, 1000 data points for corresponding images, we establish the co-occurrence matrix for 7 user intent classes defined in Section 3.3. In this aspect, the co-occurrence matrix and the associated Markov chain are as shown in Table 1 and Table 2, respectively.

By using an iterative method, we obtain the stationary distribution as:

$$\pi = [\pi_1, \pi_2, \dots, \pi_7] = [0.1505 \quad 0.2035 \quad 0.1975 \quad 0.1460 \quad 0.09647 \quad 0.1126 \quad 0.0934]$$

Figure 3 shows that the Query Intent for Information Search is represented by Q1, Q2, Q3, Q5 and Q7. Thus the percentage of Query Intent for Informational Class = $(Q1 + Q2 + Q3 + Q5 + Q7) = 0.741414 = 74\%$. Similarly we have,

The percentage of Query Intent for Navigational Class = $(Q4) = 0.14596 = 15\%$

The percentage of Query Intent for Transaction Class = $(Q6) = 0.112626 = 11\%$

5. Conclusion. In this paper we had proposed a Markov chain approach to user intent query classification scheme for web image search engine. It is important to analyze image queries from image search engines to understand how users search for images on the Web and what kind of information behind the queries a user expects. We have considered that the intent of a user query can be classified into three categories: informational, navigational and transactional on basic level and further they can derive seven more categories. We have tested the query representation by using Markov stationary distribution. We

found more favorable results for informational class than navigational and transactional categories. In future, we would like to investigate the user intent query by using more than three categories.

Acknowledgment. This work is partially supported by KAKENHI 25330133 Grant-in-Aid for Scientific Research(C).

REFERENCES

- [1] C. Thao and E. V. Munson, A relevance model for web image search, *Proc. of International Workshop on Web Document Analysis*, UK, pp.57-60, 2003.
- [2] L. H. Armitage and P. G. Enser, Analysis of user need in image archives, *Journal of Information Science*, vol.23, pp.287-299, 1997.
- [3] A. Broder, *A Taxonomy of WIGIR Forum*, vol.36, pp.3-10, 2002.
- [4] B. J. Jansen and U. Pooch, A review of web searching studies and a framework for future research, *Journal of the American Society for Information Science*, vol.52, no.3, pp.235-246, 2001.
- [5] B. J. Jansen and A. Spink, How are we searching the World Wide Web? A comparison of nine search engine transaction logs, *Information Processing & Management*, vol.42, no.1, pp.248-263, 2006.
- [6] B. J. Jansen, A. Spink and T. Saracevic, Real life, real users, and real needs, a study and analysis of user queries on the web, *Information Processing & Management*, vol.36, no.2, pp.207-227, 2000.
- [7] U. Lee, Z. Liu and J. Cho, Automatic identification of user goals in web search, *Proc. of World Wide Web Conference*, Chiba, Japan, pp.391-401, 2005.
- [8] S. Hastings, Query categories in a study of intellectual access to digitized art images, *Proc. of the ASIS 58th Annual Meeting*, vol.32, pp.3-8, 1995.
- [9] M. Markkula and E. Sormunen, End-user searching challenges indexing practices in the digital newspaper photo archive, *Information Retrieval*, vol.1, pp.259-285, 2000.
- [10] D. E. Rose and D. Levinson, Understanding user goals in web search, *Proc. of the World Wide Web Conference*, NY, USA, pp.13-19, 2004.
- [11] M. Lux, C. Kofler and O. Marques, A classification scheme for user intentions in image search, *Proc. of the 28th International Conference on Human Factors in Computing Systems, Extended Abstracts*, Atlanta, Georgia, USA, pp.3913-3918, 2010.
- [12] D. Cai, X. He, Z. Li, W.-Y. Ma and J.-R. Wen, Hierarchical clustering of WWW image search results using visual, textual and link information, *Proc. of the 12th ACM International Conf. on Multimedia*, NY, USA, pp.952-959, 2004.
- [13] X. O. Tang, K. Liu, J. Y. Cui, F. Wen and X. G. Wang, IntentSearch: Capturing user intention for one-click Internet image search, *IEEE Trans. PAMI*, vol.34, no.7, pp.1342-1353, 2012.
- [14] J. Cui, F. Wen and X. Tang, Real time google and live image search re-ranking, *Proc. of the 16th ACM International Conf. on Multimedia*, Vancouver, British Columbia, Canada, pp.729-732, 2008.
- [15] J. Huang and E. N. Efthimiadis, Analyzing and evaluating query reformulation strategies in web search logs, *Proc. of the 18th ACM Conf. on Information and Knowledge Management*, NY, USA, pp.77-86, 2009.
- [16] H.-T. Pu, A comparative analysis of web image and textual queries, *Online Information Review*, vol.29, no.5, pp.457-467, 2005.
- [17] B. C. Vattikonda et al., Interpreting advertiser intent in sponsored search, *KDD'15*, Sydney, NSW, Australia, pp.2177-2185, 2015.
- [18] X. He et al., Practical lessons from predicting clicks on ads at Facebook, *Proc. of the 8th International Workshop on Data Mining for Online Advertising*, pp.1-9, 2014.