



**HAL**  
open science

# Natural language processing for usage based indexing of web resources

Anne Boyer, Armelle Brun

► **To cite this version:**

Anne Boyer, Armelle Brun. Natural language processing for usage based indexing of web resources. 29th European Conference on Information Retrieval - ECIR'07, Fondazione Ugo Bordoni; BCS-IRSG; ACM SIGIR, Apr 2007, Rome, Italy. pp.517-524, 10.1007/978-3-540-71496-5\_46 . inria-00172231

**HAL Id: inria-00172231**

**<https://inria.hal.science/inria-00172231>**

Submitted on 14 Sep 2007

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Natural language processing for usage based indexing of web resources

Anne Boyer and Armelle Brun

INRIA Lorraine - BP 239 - 54506 Vandœuvre lès Nancy, France  
{boyer, brun}@loria.fr

**Abstract.** The identification of reliable and interesting items on Internet becomes more and more difficult and time consuming. This paper is a position paper describing our intended work in the framework of multimedia information retrieval by browsing techniques within web navigation. It relies on a usage-based indexing of resources: we ignore the nature, the content and the structure of resources. We describe a new approach taking advantage of the similarity between statistical modeling of language and document retrieval systems. A syntax of usage is computed that designs a Statistical Grammar of Usage (SGU). A SGU enables resources classification to perform a personalized navigation assistant tool. It relies both on collaborative filtering to compute virtual communities of users and classical statistical language models. The resulting SGU is a community dependent SGU.

## 1 Introduction

The amount of available information has exponentially increased in the last years due to the development of information and communication technologies and the success all over the world of Internet applications. Therefore the identification of reliable and interesting items becomes more and more difficult and time consuming, even for skilled people using dedicated tools, such as powerful search engines. Due to the huge amount of online resources, the major difficulty is nevermore to know if a pertinent document is available but to identify the more reliable and interesting items among the overwhelming stream of available information. A key factor of success in information retrieval and delivery is the development of powerful tools easy-to-use for a large audience.

Different approaches for resources retrieval use to be explored, such as content analysis, keywords indexing and identification, topic detection, etc. [1]. A major difficulty inherent to such approaches is that one keyword may have different meanings, or not, dependent of the user, his/her context and the history of his/her past navigations. Moreover two different keywords may have similar meanings, depending on the context. Expressing a query is a difficult task for many people and a lot of research and industrial projects deal with query assistance. Furthermore automatic indexing of multimedia resources is still a hard research problem. To cope with these difficulties (query expression, multimedia indexing, etc.) we decide to investigate another way by ignoring the content, the nature, the format and the structure of resources.

This paper describes our intended work, relying on our past researches both on collaborative filtering [2, 3] and statistical language modeling [4, 5]. Our objective aims at providing a new web browsing tool based on an analysis of usage. This tool enables multimedia information retrieval by browsing techniques without expressing any query. It means that users are modeled without requiring any preferences elicitation. This approach enables to easily manage heterogeneous items (video, audio, textual, multimedia) with a single treatment, this is an advantage since classical methods require dedicated tools for resource tagging.

We plan to extract frequent patterns of consultations by taking advantage of the analogy between language-based statistical modeling and resource retrieval. These frequent patterns will allow the design of syntax of usage, relying on the hypothesis that there is logic and coherency defining implicit "rules" inside a navigation. The resulting Statistical Grammar of Usage (SGU) enables a classification, clustering and selection of resources to design personalized filtering.

In the next section, the problem of retrieving resources when browsing is stated and our approach based on the use of statistical language models is detailed. The following section presents the most popular statistical language models and their appropriateness to web browsing. Section 4 puts forward the community-based Statistical Grammar of Usage we design. Then, discussion and perspectives conclude the paper.

## 2 Principle of our web browsing tool

Our web browsing tool helps users during a navigation process: it suggests the pertinent items to a specific user, given his/her past navigation and his/her context. The aim is to compute the **pertinence** of any resource. The pertinence of a resource is defined as the interest of a user for it and allows to compute **predictions** of resources (the highest the pertinence of a resource is, the highest is its probability to be suggested to the user).

First, we hypothesize an **implicit search**, it means that the active user has no explicit queries to formulate. Secondly, we consider as a **consultation** the sequence of one or more items, dedicated to a given search. A multi-navigation is the mix of different consultations within a single browsing process. A **resource** is any item (textual, audio, video or multimedia document, web page, hyperlink, forum, blog, website, etc.), viewed as an elementary and indivisible entity without any information about its format, its content or any semantic or topic indexing. The only data describing *a priori* a resource is a normalized mark called **identifier**, enabling to identify and to locate it. Our approach relies on an analysis of usage. A **usage** is any data, explicitly or implicitly left by the user during navigation. For example, history of consultation, click-stream or log files are implicit data about the interest of the visited items for the active user. This measure can be either an explicit information as votes, annotations or any estimation computed from implicit data [6].

An advantage of our approach is that it only takes into account a measure of the user's interest for a given resource, which is directly linked to the pertinence

criterion: the user's satisfaction. Let us remember that we decide to ignore any structural or thematic information about a resource. Our approach computes a personalized indexing of resources not in terms of its intrinsic nature but in terms of a more subjective but more reliable and pertinent criterion, i.e. the user's context, preferences and habits. It is the reason why this approach manages heterogeneous resources with a single treatment.

The question to solve is the following: how to estimate the *a priori* pertinence of a resource for a given user. The difficulty relies on sparsity of data: we don't have any appreciation of a resource if this user has not seen it and usually many resources have not been seen by this user. To compute the *a priori* pertinence of a resource, we plan to design a grammar of usage. As a **grammar of language** is the set of rules describing the relation between words, a **grammar of usage** is the set of rules describing the relation between resources. A grammar of language estimates if a word is pertinent given the beginning of a sentence. A grammar of usage allows to estimate if a resource is relevant for a specific user given his/her previous consultations. There is no *a priori* grammar of usage, as Internet is a dynamic and moving environment. A means to cope with the difficulty of designing an *a priori* grammar is the use of a statistical approach based on usage analysis. As huge usage corpora are available (log files, clickstream, etc.) it makes it possible to explicit regularities in terms of resource consultations. This statistical approach can be investigated in a similar way to language modeling based on statistical models.

The resulting grammar is called a **Statistical Grammar of Usage** (SGU). It enables the computation of the probability of a resource given the active user and his/her sequence of navigation. This probability measures the pertinence of the resource. A SGU, if trained on the whole usage corpus, is a general grammar since it is learned for all users in all contexts. The accuracy of such a grammar is insufficient and furthermore, the presupposed logic and coherency between users becomes a too strong and unrealistic hypothesis. Given two users, it seems unlikely that they exhibit the same resource consultation behavior: the SGU has to be personalized. Nevertheless, learning a user-specific SGU requires a large amount of data for each user and it is unrealistic to wait for collecting enough data to train it. It is the reason why we will determine groups of users with similar behavior called **communities**. Thus we plan to compute a SGU for each community and design a **community-based SGU**. Users are preclassified into a set of coherent communities, in terms of resource consultation behavior. Collaborative filtering techniques are a means to build coherent communities in terms of usage. This gathering can be compared to topic classification in natural language processing.

The principle of collaborative filtering techniques [7] amounts to identifying the active user to a set of users having the same tastes and, that, based on his/her preferences and his/her past visited resources. This approach relies on a first hypothesis that users who like the same documents have the same topics of interests and on a second hypothesis that people have relatively constant likings. Thus, it is possible to predict resources likely to match user's expectations by

taking advantage of experience of his/her community.

A first comment on usual collaborative filtering techniques is that the structure of navigation is ignored. However, this aspect can be crucial in some applications such as web browsing. For example, a user may not like a resource because he/she has not previously read a prerequisite resource. Thus the SGU will submit a resource when it becomes pertinent for a user, for example when he/she has read all prerequisites. As statistical language models emphasize the order of words in sentences, it seems interesting to determine if such models and collaborative filtering can be used together to improve the quality of suggestions.

### 3 Statistical language models

#### 3.1 Overview of statistical language models

The role of a statistical language model (**SLM**) is to assign a likelihood to a given sentence (or sequence of words) in a language [8]. A SLM is defined as a set of probabilities associated to sequences of words. These probabilities reflect the likelihood of those sequences. SLM are widely used in various natural language applications such as automatic translation, automatic speech recognition, etc. Let the word sequence  $W = w_1, \dots, w_S$ . The probability of  $W$  is computed as the product of the conditional probabilities of each word  $w_i$  in the sentence. To estimate these probabilities, a vocabulary  $V = \{w_j\}$  is stated. The probability of the sequences of words are trained on a training text corpus.

#### 3.2 How can web browsing take advantage of SLM ?

Web browsing and statistical language modeling domains seem to be similar in several points. First, statistical language modeling uses a vocabulary made up of words. This set can be viewed as similar to the set of resources  $R$  of the web. Then, the text corpus is made up of sentences of words, they can be viewed as similar to the sequences of consultations of the usage corpus. A sequence of  $S$  words in a sentence is similar to a sequence of consultation of  $S$  resources. Finally, the presence of a word in a sentence mainly depends on its previous words, as the consultation of a resource mainly depends of the preceding consultations. Given these similarities, we can naturally investigate the exploitation of SLM into a web browsing assistant. As noticed in the previous section, these models have the characteristic that the order of the elements in the history is crucial. This aspect may be important for specific resources in web browsing.

However, we have to notice that web browsing and natural language processing have two major differences. The first one is that it is possible that a user may mix different queries within a single history (we will call this "multi-navigation") but it is unrealistic to mix different sentences when speaking or writing. This first remark brings us to consider a generalization of SLM to integrate "multi-navigation" in the browsing process. The second one is that natural language exhibits strongest constraints: each word in a sentence is important and deleting

or adding a word may change the meaning of the sentence. Web browsing is not so sensitive and adding or deleting a specific resource within a navigation may have no impact. Then we will have to consider permissive models, able to take into account less constrained histories such as navigation has.

### 3.3 n-grams language models

Due to computational constraints and probability reliance, the whole history  $h_i$  of  $w_i$  cannot be systematically used to compute the probability of  $W$ . Classical SLM aim at reducing the size of the history while not decreasing performance.

$n$ -grams models reduce the history of words to their  $n - 1$  previous words. These models are the most commonly used in most of natural language applications.  $n$ -grams model can be directly used in web browsing assistance. In the previous section, we put forward that the quality of the model will be increased if it is dedicated to a community and trained on the corresponding community usage corpus. Thus, the usage corpus is split into community usage corpora and a model is trained on each community corpus.

Let a community  $c_j$  and a sequence of consultations of resources  $h_j = R_{j1}, \dots, R_{ji-1}$ . For each resource  $R_i \in R$ , the  $n$ -grams model computes the probability  $P_n(R_i | R_{i-n+1}, \dots, R_{i-1}, c_j)$ . The history  $h_j$  is reduced to the  $n - 1$  last resources consulted, other resources are discarded. Thus, this model assumes that the consultation of a resource  $R_i$  does not mainly depend on resources consulted far from  $R_i$ .

As previously mentioned, the behavior of users is less constrained than language: adding or deleting a resource in a sequence of consultations has a lower influence on the result of the search than adding or deleting a word in a sentence. This model does not ideally match our retrieval problem since the history considered is the exact sequence of consultations  $R_{i-n+1} \dots R_{i-1}$ , that may be too restrictive in the general case. However, this model may be suitable for frequent sequences of consultations, that can be considered as “patterns of consultation”. They are assigned a high probability, thus increasing the probability of resources inside such sequences. It should be interesting to take into account, in a more adequate way, such “patterns of consultations”.

As  $n$ -grams models exhibit strong constraints, we are also interested in more permissive models. Trigger-based language models seem to me more adequate to less constraint histories such as navigation.

### 3.4 Trigger-based language models

Trigger-based models [9] aim at considering long-time dependence between two words ( $w_x$  and  $w_y$  for instance). Dependence is measured by Mutual Information (MI) [10]. This measure can easily integrate long-time dependence by using a distance parameter  $d$ .  $d$  is the maximum number of words occurring between  $w_x$  and  $w_y$ , a window of  $d$  words is thus considered.

A couple  $(w_x, w_y)$  with a high MI value means that  $w_x$  and  $w_y$  are highly correlated and the presence of  $w_x$  raises the probability of occurrence of  $w_y$ , at a

maximal distance of  $d$  words.  $(w_x, w_y)$  is named a trigger. This model considers only highly correlated pairs of words (corresponding to high MI values), useless pairs are discarded. The resulting set is called  $S$ .

Given history  $h_j = w_1, \dots, w_{i-1}$ , the trigger model computes the probability of  $w_i$  as:

$$P_t(w_i | h_j) = \frac{\sum_{w_j \in h_j} \delta_{w_j, w_i, h_j, S}}{\sum_{w_j \in h_j} \sum_{w_t \in V} \delta_{w_j, w_t, h_j, S}} \quad (1)$$

$$\text{with } \delta_{w_j, w_i, h_j, S} = \begin{cases} 1 & (w_j, w_i) \in S \text{ and } d_j(w_j, w_i) \leq d \\ 0 & \text{otherwise} \end{cases}$$

where  $d_j(w_j, w_i)$  is the distance between  $w_j$  and  $w_i$ , in terms of words in  $h_j$ .

In our web browsing assistant tool, the trigger model is made up of triggers of resources  $(R_x, R_y)$ . The consultation of  $R_x$  triggers the consultation of  $R_y$ , at a maximal distance of  $d$  resources. As MI measure is not symmetric ( $MI(R_x; R_y) \neq MI(R_y; R_x)$ ), this model integrates order between resources, that may be crucial for specific resources.

The advantage of such a model is the long-time dependence between both resources. In a consultation, two resources can be viewed with various values of distance without changing the meaning of the consultation. Trigger models enable to modelize this kind of influence, when the value of the distance between items is not discriminant but the order of occurrence is meaningful. Such a model is less constrained than  $n$ -grams models and seems to be adequate to the navigation problem.

Similarly to  $n$ -grams model, a trigger-model is developed for each community  $c_j$ . MI values are computed for each couple of resources and for each community. A set of the most related triggers ( $S_{c_j}$ ) is extracted for each community  $c_j$ .

The probability of a resource  $R_i$ , given the community  $c_j$ , its corresponding set of triggers  $S_{c_j}$  and the sequence of consultation of resources  $h_j = R_1, \dots, R_{i-1}$  is:

$$P_t(R_i | h_j, c_j) = \frac{\sum_{R_x \in h_j} \delta_{R_x, R_i, h_j, S_{c_j}}}{\sum_{R_x \in h_j} \sum_{R_y \in R} \delta_{R_x, R_y, h_j, S_{c_j}}} \quad (2)$$

$$\text{with } \delta_{R_x, R_i, h_j, S_{c_j}} = \begin{cases} 1 & (R_x, R_i) \in S_{c_j} \text{ and } d_j(R_x, R_i) \leq d \\ 0 & \text{otherwise} \end{cases}$$

where  $d_j(R_x, R_i)$  is the distance between  $R_x$  and  $R_i$  in history  $h_j$ .

## 4 Towards a community-based SGU

The SGU we propose in this article has the advantage of considering both the community of the active user and his/her consultation history (sequence of consultation), whereas state of the art models usually exploit the set of consultations. The use of this model relies on two steps:

#### 4.1 Determination of the community of the active user

The first objective is to compute a set of user communities based on an analysis of usage. To achieve this goal, we use collaborative filtering techniques. The set of users is split into classes by using a recursive k-means like algorithm [2], the similarity between two users is estimated as the mean of the distance for each commonly voted resource [11]. The whole corpus is then split into community sub-corpora. Each one is made up of usage of any user of the community. A user is then assigned to the closest community using the same similarity measure.

#### 4.2 Computation of the probability of a resource

Given the community  $c_j$  of user  $U_j$ , and his history  $h_j$ , the computation of the probability of a resource  $R_i$  relies on three sub-models based on language models presented in section 4. The first sub-model computes the probability  $P_n(R_i | h_j, c_j)$ , by exploiting the probabilities of resources sequences of the  $n$ -grams model. The second sub-model is the trigger model, it computes the probability  $P_t(R_i | h_j, c_j)$ . The last sub-model is devoted to resources out of the training corpus. A probability *a priori*  $P_a(R_i | c_j)$  is set to each resource  $R_i \in R$ . The resulting model, that can be viewed as a community-based SGU, computes the linear combination of the three previously described sub-models.

$$P(R_i | h_j, c_j) = \lambda_n P_n(R_i | h_j, c_j) + \lambda_t P_t(R_i | h_j, c_j) + \lambda_a P_a(R_i | c_j) \quad (3)$$

where  $\lambda$  are optimized with EM algorithm [12] on a development corpus.

Thus, given a user  $U_j$  and his/her history  $h_j$ , we first have to determine the community  $c_j$  he/she belongs to. Then, the probability of any available resource is computed given the SGU learned for this community.

Then, the  $N$  most likely resources are selected. The systematic selection of resources in the same subset of likely items avoids the introduction of novelty in resources suggestion. To enable novelty in suggestion, we randomly select a subset of unlikely resources (SUR) that is added to the previous subset of  $N$  likely resources (SLR) to build the set of candidates (SC). We determine the suggested resources for a specific user using a roulette wheel. We assign a sector of the wheel to any resource in SC; the size of this sector is proportional to the probability of occurrence of this resource as given by the SGU. One or several resources from SC are then drawn independently using this roulette wheel principle and are submitted to user  $U_j$ .

## 5 Discussion and perspectives

This paper aims at describing a new web browsing assistant, based on usage and natural language processing. This approach exempts the difficult task of content indexing and facilitates heterogeneous resources management. Similarities between SLM and web browsing are put forward, therefore the integration of SLM is investigated. The resulting model is a Statistical Grammar of Usage (SGU).



As a single SGU may be unefficient, it has to be personalized. To tackle sparsity of data, a preclassification of users into communities is performed. Community-based SGU are then proposed. A second contribution consists in the design of community-based SGU, predicting the sequentiality of resources during navigation. Moreover, a community-based SGU builds an *a posteriori* structure of navigation based on the subjective but reliable measure of pertinence of a resource for a user. Consequently it performs a personalized indexing of resources, based on usage analysis.

As collaborative filtering techniques used to build communities and triggers used to suggest resources have both proved their efficiency in their respective domain, a first perspective is the validation of the community-based SGU in terms of quality of predictions in web browsing. A second perspective is the use of the community-based SGU to compute a personalized classification of resources, depending not only on topics but also on user's preferences and context.

## References

1. R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*, ACM Press, New York, 1999.
2. S. Castagnos and A. Boyer, "A client/server user-based collaborative filtering algorithm model and implementation," in *Proceedings of the 17th European Conference on Artificial Intelligence (ECAI 2006)*, Riva del Garda, Italy, august 2006.
3. S. Castagnos and A. Boyer, "Frac+: A distributed collaborative filtering model for client/server architectures," in *2nd conference on web information systems and technologies (WEBIST 2006)*, Setbal, Portugal, 2006.
4. K. Smaïli, A. Brun, I. Zitouni, and J.P. Haton, "Automatic and manual clustering for large vocabulary speech recognition: A comparative study," in *European Conference on Speech Communication and Technology*, Budapest, Hungary, 1999.
5. A. Brun, K. Smaïli, and J.P. Haton, "Contribution to topic identification by using word similarity," in *International Conference on Spoken Language Processing (ICSLP2002)*, 2002.
6. P. Chan, "A non-invasive learning approach to building web user profiles," in *5th International Conference on Knowledge Discovery and Data Mining - Workshop on Web Usage Analysis and User Profiling*, San Diego, USA, august 1999.
7. J. Herlocker, J. Konstan, L. Terveen, and J. Riedl, "Evaluating collaborative filtering recommender systems," *ACM Transactions on Information Systems (TOIS)*, vol. 22, no. 1, pp. 5–53, january 2004.
8. R. Rosenfeld, "Two decades of statistical language modeling: Where do we go from here," 2000.
9. R. Rosenfeld, "A maximum entropy approach to adaptative statistical language modeling," *Computer Speech and Language*, vol. 10, pp. 187–228, 1996.
10. N. Abramson, *Information Theory and Coding*, McGraw-Hill, New-York, 1963.
11. U. Shardanand and P. Maes, "Social information filtering: algorithms for automating "word of mouth"," in *Proceedings of the ACM CHI'95 - Conference on Human Factors in Computing Systems*, 1995, vol. 1, pp. 210–217.
12. A.P. Dempster, N.M. Laird, and D.B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *J. of the Royal Statistical Society*, vol. 39, 1977.

This article was processed using the L<sup>A</sup>T<sub>E</sub>X macro package with LLNCS style