

learning” (eg. entertainment recommendation). Holte and Drummond [15], Drummond et al. [9] designed a system that assists browsing of software libraries, taking keywords from the user and using a rule-based system with forward chaining inference, assuming that the library consists of one type of item and the user goal is a single item. Etzioni and Weld [11] offer an integrated interface to the Internet combining UNIX shell and the World Wide Web to interact with Internet resources. Their agent accepts high-level user goals and dynamically synthesizes the appropriate sequence of Internet commands to satisfy those goals. Hammond et al. [?] and Burke et al. [6] developed a system that uses a “natural language question-based interface to access distributed text information sources” and helps the user to find answers to her/his question in a databases such as FAQ files.

One of the available information sources is the World Wide Web and it is currently growing quickly, attracting many users with different interests. Since interaction with the World Wide Web (WWW) is by means of computer, one can use computers to “observe” and record user actions and use this information to help users find their way in the WWW.

Personal WebWatcher is a system that observes users of the WWW and suggests pages they might be interested in. It learns user interests from the pages requested by the user. The learned model of user interests is then used to suggest hyperlinks on new HTML-pages requested by and presented to the user via Web browser that enables connection to “proxy” eg. NETSCAPE. Section 2 gives an overview of Personal WebWatcher, its functionality and some directions for further development. Section 3 describes some of the research problems related to Personal WebWatcher’s learning. The structure of the system and its implementations are described in Section 4 and Perl code is given in the Appendix. Results of the first experiments are given in Section 5.

2 “Personalizing” WebWatcher

Personal WebWatcher is mainly inspired by WebWatcher: “a Learning Apprentice for the World Wide Web” [2], [16] and some other work related to learning apprentice and learning from ~~text~~ [17], [19], [21], [25]. The idea of a learning apprentice is to automatically customize to individual users, using each user interaction as a training example.

WebWatcher can be described as an agent that assists users in locating information on the WWW. It learns by observing a user on her/his way through the WWW and suggests interesting hyperlinks whenever it is confident enough. The idea is that the user provides a few keywords describing a search goal and WebWatcher highlights related hyperlinks on the current page and/or adds new hyperlinks to the current page. It can also suggest pages related to the current page using information stored in the structure of hypertext without considering the text itself [16], or send an e-mail message to the user whenever specified pages change. The same WebWatcher version is designed to serve all users, collecting information and sharing it between users. For example, if someone recognizes a page as being related to the keywords she/he typed in the system at the beginning of the search, this can be useful for any user searching for similar

information.

Unlike WebWatcher, Personal WebWatcher (PWW) is structured to specialize for a particular user, modeling her/his interests. It “watches over the user’s shoulder” the similar way WebWatcher does, but it avoids involving the user in its learning process (it doesn’t ask the user for any keywords or opinions about pages). It solely records the addresses of pages requested by the user and highlights hyperlinks that it believes will be of interest. In the learning phase (typically during the night), requested pages are analyzed and a model of user interests is generated/updated. This model is used to give advice for hyperlinks on retrieved HTML-pages requested by and presented to the user via Web browser.

Since each user has her/his own copy of the system - her/his own agent, these agents can communicate and exchange information on a base of similarity between their users, often referred to as collaborative or social learning [22], [31]. There are also some other types of “Personal agents” the user could use, for example, an agent for Calendar Apprentice [25], and these agents can exchange information they have about the same user in different fields/activities. we plan to investigate these in the future work on Personal WebWatcher.

3 PWW “Behind the stage”

There are many research question behind Personal WebWatcher, that wait to be answered. Here, we consider some of them that are related to learning. There are also many others, for example, how to design communication between user and agent and between different agents, how to provide privacy to the user, to which extent agent should involve user to its learning/exploration process, ...

If we concentrate on learning, we first want to know what kind of problem are we faced with, how to represent it to some learning algorithm and which kind of algorithm to use. Currently we restricted our work to text documents (plain text and HTML) so we are faced with the problem of text-learning having mainly short to medium size documents with varying vocabulary. This section gives different approaches related to (1) document representation, (2) feature selection and (3) learning used on text-learning problem. Table 1 summarizes them over some related papers in order to give an idea about current trends.

3.1 Document representation

The frequently used document representation in Information Retrieval and text-learning is the so called *TFIDF*-vector representation. It is a *bag-of-words representation*: all words from the document are taken and no ordering of words or any structure of text is used (see Figure 1). Since most of our documents are in HTML format, there is a well defined underlying structure that could be used. There is also additional information in plain text, for example structure of sentences, position of words or neighboring words. The question is how much can we gain considering additional information in learning (and what information to consider) and what is the price we have to pay for it? There is currently no well studied comparison or directions that we are aware of. There is

| Paper reference | Doc. Representation | Feature Selection | Learning |
|---|---|---|--|
| Apté et al. [1] | bag-of-words (frequency) | stop list+ frequency weight | Decision Rules |
| Armstrong et al. [2] | bag-of-words (Boolean) | mutual info. | TFIDF, Winnow, WordStat |
| Balabanovic et al. [3] | bag-of-words (frequency) | stop list+stemming+ keep 10 best words | TFIDF |
| Bartell et al. [4] | bag-of-words (frequency) | latent semantic indexing using SVD | — |
| Berry et al. [5] Foltz and Dumais [12] | bag-of-words(frequency) | latent semantic indexing using SVD | TFIDF |
| Cohen [8] | bag-of-words (Boolean) | infrequent words pruned | Decision Rules ILP |
| Joachims [17] | bag-of-words (frequency) | in/frequent words+ mutual info. | TFIDF, PrTFIDF, Naive Bayes |
| Lewis et al. [23] | bag-of-words (Boolean) | log likelihood ratio | logistic regression combined with Naive Bayes |
| Maes [24] | bag-of-words+ header info. | mail/news header info.+ selecting keywords | Memory-Based reasoning |
| Pazzani et al. [26] | bag-of-words (Boolean) | stop list+ mutual info. | TFIDF, Naive Bayes, Nearest Neighbor, Decision Trees |
| Sorensen and Mc Elligott [33], [10] | n-gram graph (only bigrams) | weighting graph edges | connectionist combined with Genetic Algorithms |
| Yang [35] | bag-of-words (Boolean, frequency, TFIDF) | stop list | adapted k-Nearest Neighbor |

Table 1: Document representation, feature selection and learning algorithms used in some related work.

some evidence in Information Retrieval research, that for long documents, considering information additional to bag-of-words is not worth efforts. But our documents are mostly HTML-documents on the WWW and they are not especially long!

In the process of using text to learn how to give advice for a hyperlink, different approaches can be used to decide which part of text to use and how to represent it. Personal WebWatcher uses an approach similar to that of WebWatcher. In WebWatcher the bag-of-words representation is used, where considered text consists of underlined words, words in the sentence containing the hyperlink, words in all the headings above the hyperlink and words given as keywords by the user [2]. Some later versions of the WebWatcher system change slightly the way of constructing text for learning, eg. adding words in the document retrieved behind hyperlink. Many current systems that learn on text use the bag-of-words representation using either Boolean features indicating if specific word occurred in document (eg. [2], [8], [23], [26], [35]) or frequency of word in a given document (eg. [1], [3] [4], [5], [17], [35]). There is also some work that uses additional information such as word position [8] or word tuples called n-grams [33].

We decided to use the bag-of-words representation using frequency of word and observe success of given advice (whether user selected the advised hyperlink). In case of poor system performance, some additional information from HTML-structure could be added, for example, frequency of word in headlines of a given document. We would

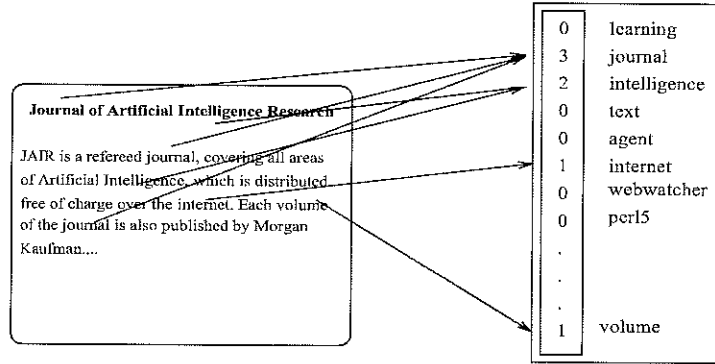


Figure 1: Bag-of-words representation using frequency vector.

like to spend time on extracting and using this additional information only when highly probable that we will gain a decent improvement in system performance. This is currently under research.

3.2 Feature selection

If we decide to ignore all the additional information and use the bag-of-words approach, we still end up with several tens of thousands of different words that occur in our documents. Not only is using all these words time-consuming but also many of them are not really important for our learning task. One of the approaches to reduce number of different words is to use “stop-list” containing common English words (eg. a, the, with) or pruning the infrequent and/or very frequent words ([8], [17]). There is also the possibility of word stemming. Many approaches introduce some sort of word weighting and select only the best words (eg. [1], [2], [3], [17], [23], [26]) or reduce dimensionality using latent semantic indexing with singular value decomposition (eg. [4], [5], [12]). Some of the Machine Learning techniques for feature selection could also be used (eg. [7], [32], [18]) but most of them take too long in situations with several tens of thousands of features.

In the current implementation of PWW we weight words using mutual information between word occurrence and class value [34], the same way as used in [2] or [17]. The research topic is the number of best words to consider and we observe its influence on classification accuracy and precision of the best suggested hyperlink (see Section 5). Mutual information assigns higher weight to the words that make better distinction between interesting and uninteresting documents. One of the practical problems during classification is that a new document often contains very few or even none of the selected words. Since we are more interested in positive class (interesting documents) and we want to have words that are frequent, it might be better to include in the weighting formula the probability of a word occurring in the positive class or frequency of the word ([23]) $TF(w) \log(\frac{P(w|c)}{P(w|\bar{c})})$, where w is a selected word, c is the positive class and $TF(w)$ is the frequency of word w .

We plan to make additional experiments using the proposed combined weighting as

well as using some other weighting methods. For example, combining a stop list with weighting words by their frequency and keeping the most frequent words [35] or using word weighting used in the odds ratio method [27]. Odds ratio is the method of document ranking according to their relevance for a given problem (eg. being interesting for user). Ranking of documents is defined as $ranking(d, c) = \log \frac{P(c|d)}{P(\bar{c}|\bar{d})} = \log \frac{P(c)P(d|c)}{P(\bar{c})P(\bar{d}|\bar{c})} = \log \frac{P(c)\prod_{w \in d} P(w|c)}{P(\bar{c})\prod_{w \in \bar{d}} P(w|\bar{c})} = \sum_{w \in d} weight(w) + const.$ Word weight that it defines can be used for feature selection, weighting all words and selecting words with the highest weight. This word weight is defined as

$$weight(w) = \log \frac{p(w)(1 - \bar{p}(w))}{(1 - p(w))\bar{p}(w)}$$

where $p(w) = \frac{TF(w, c) + const.}{\#docs(c) + 1}$, where $docs(c)$ is number of documents in class c and $\bar{p}(w)$ is the same as $p(w)$ except that c is substituted with \bar{c} . Shaw [28] proposed special handling of singularities in the above formulas for $p(w)$ and $\bar{p}(w)$, namely, $p(w) = \frac{1}{\#docs^2}$ when $TF(w, c) = 0$ and $p(w) = 1 - \frac{1}{\#docs^2}$ when $TF(w, c) = \#docs(c)$

3.3 Learning algorithm

One of the well-established techniques for text in Information Retrieval is to represent each document as a *TFIDF*-vector in the space of words that appeared in training documents [30], sum all interesting document vectors and use the resulting vector as a model for classification (based on [29] relevance feedback method). Each component of a document vector $d^{(i)} = TF(w_i, d)IDF(w_i)$ is calculated as the product of *TF* (Term Frequency — number of times word w_i occurred in document) and *IDF* = $\log \frac{D}{DF(w_i)}$ (Inverse document Frequency), where D is the number of documents and document frequency $DF(w_i)$ is the number of documents word w_i occurred in at least once. The exact formulas used in different approaches may slightly vary (some factors are added, normalization performed [30]) but the idea remains the same. A new document is then represented as a vector in the same vector space as the generated model and the distance between them is measured (usually defined as a cosine of angle between vectors) in order to classify the document. This technique has already been used in Machine Learning experiments on World Wide Web data (eg. [2], [3], [5], [17], [26]).

There are also some other techniques for model generation that have been used in text-learning. Armstrong et al. [2] used a statistical approach they called Word-Stat that assumes mutual independence of words and defines probability of class c as $P(c) = 1 - \prod_w (1 - P(c/w))$. Pazzani et al. [26] used a Naive (Simple) Bayesian classifier on Boolean vectors, that assumes independence of words and defines probability of class c for given document doc that contains words w as proportional to $P(c)\prod_w P(c/w)$. They also used Nearest Neighbor and symbolic learning using Decision Trees. A variant of k-Nearest Neighbor was also used by Yang [35], where relevance of class c given document doc is defined as $rel(c/doc) = \sum_{i=1}^k similarity(doc, D_i)P(c/D_i)$ and *similarity* is measured by cosine between vectors and $P(c/D_i) = \frac{\#D_i Inc}{\#D_i In training data}$ (same document may occur several times being classified in different categories). Joachims [17] introduced Probabilistic *TFIDF* that takes into account document representation Θ and de-

finds probability of class c for given document doc that contains words w as $P(c/doc) = \sum_w \frac{P(c)P(w/c)}{\sum_i P(c_i)P(w/c_i)} P(w/doc, \Theta)$ where $P(w/c) = \frac{TF(w,c)}{\sum_i TF(w_i,c)}$ and $P(w/doc, \Theta) = \frac{TF(w,doc)}{\sum_i TF(w_i,doc)}$. He also used Naive (Simple) Bayesian classifier on frequency vectors, the same as we used in PWW. It assumes independence of words and defines probability of class c for given document doc that contains words w as

$$P(c/doc) = \frac{P(c)\prod_w P(w/c)^{TF(w,doc)}}{\sum_i P(c_i)\prod_w P(w/c_i)^{TF(w,doc)}}$$

where $P(w/c) = \frac{1+TF(w,c)}{\#words + \sum_i TF(w_i,c)}$; $TF(w,c)$ denotes frequency of word w in documents of class c and $TF(w,doc)$ denotes frequency of word w in document doc .

Apté et al. [1] used Decision Rules and observed that in case of different topics being categories, it is better to select features for each given topic (using stop-list and frequency weighting) than for all topics at once, even if the set of features is additionally reduced for each topic using entropy-based measure to weight features. Cohen [8] used Decision Rules and the Inductive Logic Programming systems FOIL and FLIPPER. Lewis and Gale [23] used a combination of a Naive Bayesian classifier and logistic regression defining probability of class c for given document doc that contains words w as $P(c/doc) = \frac{\exp(a+b \sum_w \log \frac{P(w/c)}{P(w/c_i)})}{1 + \exp(a+b \sum_w \log \frac{P(w/c)}{P(w/c_i)})}$.

Maes [24] used Memory-Based reasoning, McElligot and Sorensen [10], [33] used a connectionist approach combined with Genetic Algorithms.

We decided to test different learning algorithms on PWW data (see Section 5), since it is not clear which algorithm is the most appropriate. The current version of PWW uses a Naive (Simple) Bayesian classifier on frequency vectors to generate a model of user interests, that is used for advising hyperlinks.

4 Structure of Personal WebWatcher

Personal WebWatcher consists of two main parts: a **proxy server** that interacts with the user via Web browser and a **learner** that provides the user-model to the server (see Figure 2). Communication between them is via disk; the **proxy** saves addresses of visited documents (URLs) and the **learner** uses them to generate model of user interests. The whole system is implemented in approximately 2500 lines of Perl code and 1500 lines of C++ code.

The **proxy server** consists of three main parts, each implemented as a Perl script: **proxy** (additionally calls external **fetcher** code to fetch the page), **adviser** and **classifier** (calls external C++ code for classification). **Proxy** waits in an infinite loop for a page request from browser. On request, it fetches the requested document and; if it is an HTML-document adds advice and; forwards the document to the user. To add advice **proxy** forwards the page to **adviser**, that extracts hyperlinks from document and calls external code for **classification** that uses generated user-model. A limited number of hyperlinks that are scored above some threshold are recommended to the

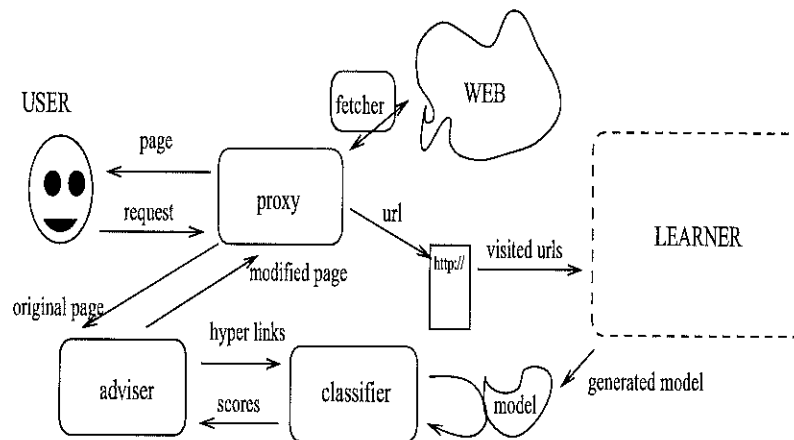


Figure 2: Structure of Personal WebWatcher. The learning part is described separately

user, indicating their scores using graphical symbols placed around each advised hyperlink. For example, in Figure 3 three hyperlinks are suggested by PWW: “Machine Learning Information Services” and two project members (Dayne Freitag, Thorsten Joachims). There is a banner at the top of the page showing that PWW is “watching over the user’s shoulder”.

4.1 Structure of the learning module

Learner works in two versions: learning a new model from scratch (**LEARNER**) or updating an existing model (**UPDATER**) as shown in Figure 4. The difference is that the first one has to define the domain (words to be used) and starts learning with an empty model, while the second one has already defined which words to use in representing documents as frequency vectors and has an existing model to modify. Both versions fetch visited documents and documents one step behind the hyperlinks of visited documents and store them as positive or negative examples of user interests, depending whether the user visited the document or not (getDoc in Figure 4). Hyperlinks from visited HTML-documents are extracted and stored in an extended hyperlink format, the same that is used by **adviser**. Each hyperlink is represented in extended format, taking into account underlined words, words in a window around a hyperlink and words in all the headings above the hyperlink. Using hyperlinks represented only with underlined words is often a bad idea, eg. click here...

Hyperlinks whose documents were visited by the user are considered to be positive examples, and all the other to be negative examples of the user interests. The idea is that all hyperlinks were presented to the user and the user chose to visit some of them that meet her/his interests. This simplification is introduced to minimize users

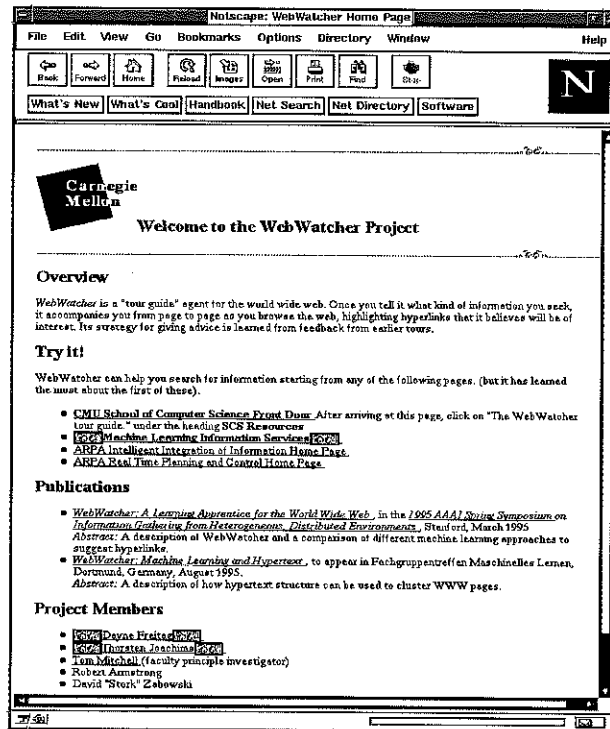


Figure 3: Example of HTML-page presented to the user by PWV.

involvement in the learning process and enable learning without asking user for page rating. We plan to make it optional in some later versions of the system.

LEARNER transforms documents into examples in two phases: (1) (docs2exs and docs2addexs in Figure 4) parsing each document, assigning an index to each word and representing it in three files as a line of word indices containing: all words, only headline words, only underlined words. (2) (exs2vec in Figure 4) calculating score (eg. information gain) for each word, selecting some top words and represent documents as bag-of-words keeping frequency for each of the top words.

UPDATER uses given top words (domain definition) and represent documents as bag-of-words, the same way LEARNER does (docs2ddexs in Figure 4).

During the model generation (GenModel in Figure 4), the system can ask for additional information about some words (stating which kind of information - functions and on which words - basic attributes). The kind of information it could ask for is specified as so called background knowledge eg. feature saying how many times a word occurred in document headlines (genAttr in Figure 4). This is currently under development.

Learning part consists basically of eight Perl scripts that call external C++ code for model generation/updating. Two scripts integrate parts of LEARNER (UPDATER) and the other six that are represented with rectangles in Figure 4 (getDocs, docs2exs, docs2addexs, exs2vec, genAttr, docs2ddexs). Two additional rectangles in Figure 4 GenModel, UpdateModel represent C++ code.

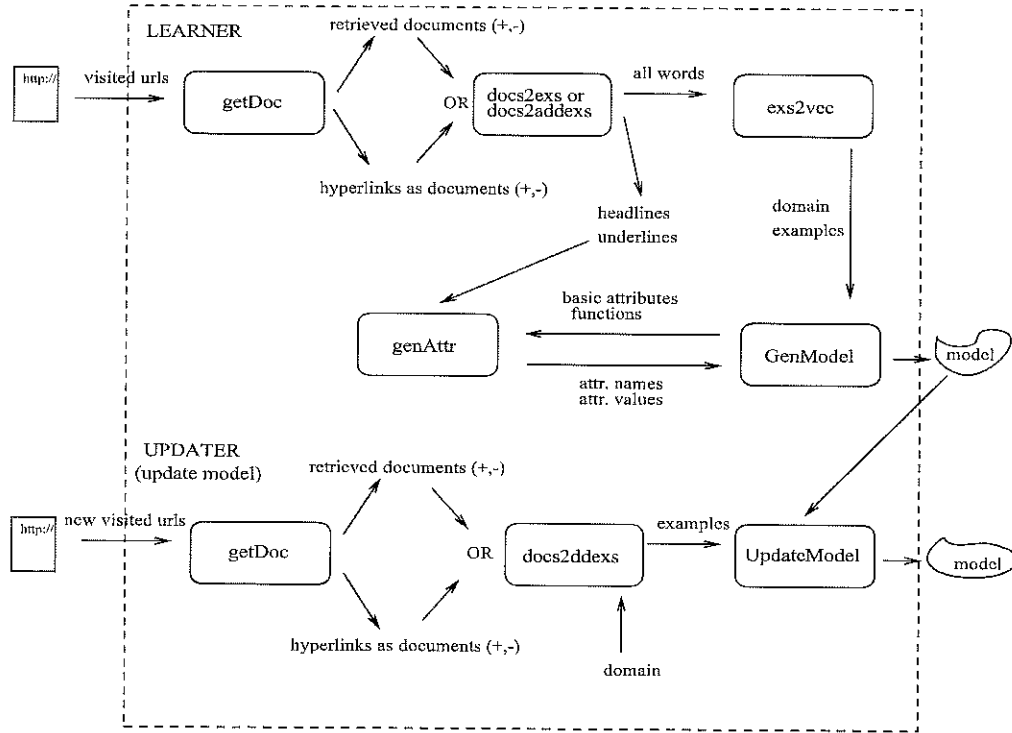


Figure 4: Structure of PWW Learner.

4.2 Model of user interests

The model of user interests is designed to predict if some document is positive or negative example of user interests. It is used to advice hyperlinks on the HTML-document requested by the user. Since the prediction should be performed while user is waiting for a HTML-document, we are actually predicting interestingness of document based on the hyperlink pointing to it, and not document itself (retrieving documents behind the requested hyperlinks is usually time consuming)

The model of user interests is generated “off-line”, usually during the night and thus its generation is not so critical in time as its usage for prediction. One of the simplest idea for learning is to use hyperlinks that occurred on the documents presented to the user as training examples and learn to predict if a new hyperlink is positive or negative example of the user interests:

$$User_{HL} : HyperLink \rightarrow \{pos, neg\}$$

What we use is an extended representation of hyperlink (see Section 4.1), that tries to capture information related to the document behind a hyperlink. But during the learning phase we can afford using more time than when adding advice, so why not retrieving documents behind hyperlinks, instead of using the extended hyperlink representation? In that case, we can learn the model of user interests directly from documents whose

interestingness we are trying to predict:

$$User_{DOC} : Document \rightarrow \{pos, neg\}$$

In this case, we end up with using the model generated from documents to predict interestingness of hyperlinks. Since our hyperlinks are represented as short HTML-documents (including headlines and some portions of text) it is not so unusual, but we could also learn a model that predicts document content based on a given hyperlink:

$$Document_{HL} : HyperLink \rightarrow Document$$

and then predict interestingness of so predicted document content using the above described model $User_{DOC}$. Our first experiments are in learning the first two models.

5 First experimental results

In order to select a good document representation, feature selection method and learning algorithm we decide to test different possibilities and compare them. Since PWW has to recommend interesting hyperlinks to the user, we are interested in measuring the precision of our system on the most highly recommended hyperlink. Precision is frequently used as a metric in Information Retrieval and it is defined as the percent of positive suggestions among all suggestions made. In our case, it is either zero (in case the best suggestion is a negative example and shouldn't be suggested to the user) or one (if we made a correct suggestion). We also measured traditional Machine Learning quality estimate classification accuracy, defined as percent of correctly classified examples (over all classes). Both quality estimates are calculated using 10-fold cross-validation technique (see Figure 5) using the same example splits for all tested algorithms.

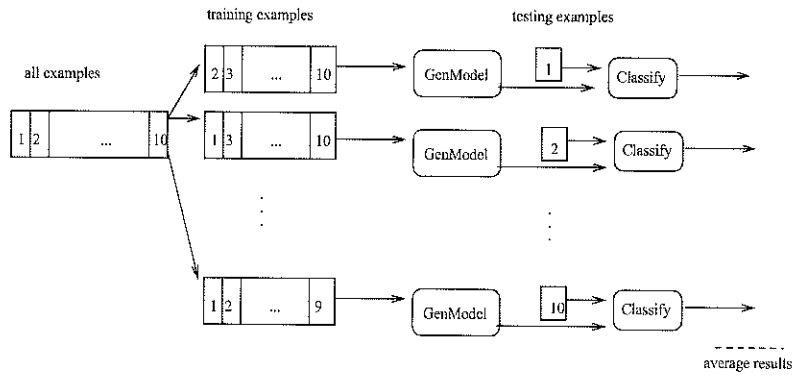


Figure 5: Illustration of 10-fold cross-validation experiments.

In first experiments we observe how vector length (number of features selected) influences model quality for two learning algorithms: Naive Bayesian classifier on frequency vector as used by Joachims [17] (see Section 3.3) and k-Nearest Neighbor approach on

frequency vectors using Euclidean distance between examples and summing class probabilities predicted by k-neighbors. We tested both algorithms on data for documents behind hyperlinks $User_{DOC}$ and data for hyperlinks $User_{HL}$ (see Section 4.2). Documents are currently represented using the bag-of-words approach (see Section 3.1) and feature selection is performed using mutual information approach (see Section 3.2). Our experiments are performed on data collected for four users participating in the HOMENET project [14] with the data characteristics given in Table 2.

| UserId and data source | probability of interestingness | number of examples | data entropy |
|------------------------|--------------------------------|--------------------|--------------|
| usr150101 | | | |
| Doc | 0.094 | 1 333 | 0.449 |
| HL | 0.104 | 2 528 | 0.480 |
| usr150202 | | | |
| Doc | 0.107 | 3 415 | 0.492 |
| HL | 0.053 | 4 798 | 0.301 |
| usr150211 | | | |
| Doc | 0.089 | 2 038 | 0.436 |
| HL | 0.044 | 2 221 | 0.259 |
| usr150502 | | | |
| Doc | 0.100 | 1 272 | 0.468 |
| HL | 0.100 | 2 498 | 0.468 |

Table 2: Data characteristics for document (Doc) and hyperlink (HL) data for each of the four HomeNet users.

In all experiments k-Nearest Neighbor achieved slightly higher classification accuracy than the Naive Bayesian classifier (see Figures 6, 7, 8, 9), but the difference is significant only in one out of six experiments (see Figure 8). Adding more than approximately 50 features doesn't appear to help for classification accuracy, but it also doesn't hurt. High classification accuracy achieved for all four users by k-Nearest Neighbor algorithm is in fact default accuracy if negative class is predicted for all documents. So what we are really interested in is making good predictions for positive documents, and that is why we decide to measure precision of the best suggested hyperlink.

Precision varies much more with vector size than classification accuracy (see Figures 10, 11, 12, 13), it actually drops in seven out of eight experiments. It seems that k-Nearest Neighbor is more stable in precision than Naive Bayesian classifier, achieving higher precision for longer vectors and about the same for up to 100 features, with exception for one user (see Figure 12), where the precision of Naive Bayesian classifier is for most vector sizes 0 for the document model and much better on short vectors than k-Nearest Neighbor for the hyperlink model.

In order to draw some conclusion about vector size and quality of algorithms, we need to perform more experiments on different users. These first experiments show that increasing vector size probably isn't as beneficial as one could expect and it even could hurt precision of the best suggestion. There is no evidence that algorithms differ substantially in classification accuracy, although k-Nearest Neighbor seems to be

more promising both in accuracy and precision. If further experiments confirm the hypothesis that long vectors are not advisable, a closer look at the short vectors (eg. see Figure 14) should give an idea about number of features that work well for both accuracy and precision.

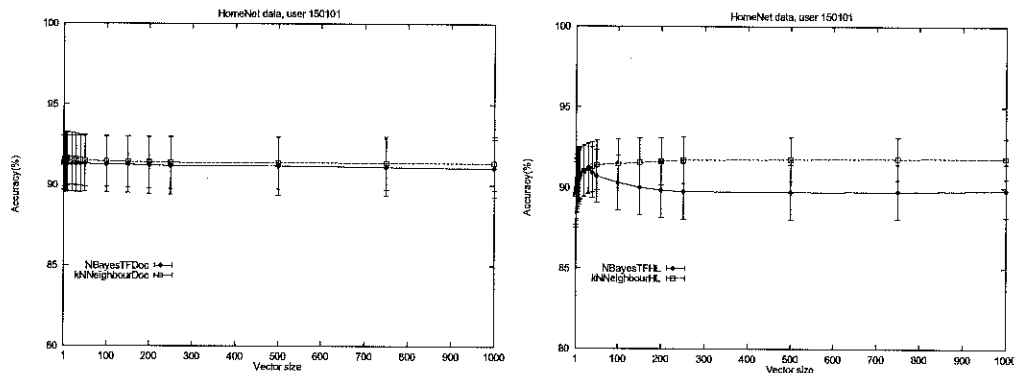


Figure 6: Influence of vector size to classification accuracy of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150101. Notice that classification accuracy scale starts at 80% accuracy. Error bars show standard deviation since accuracy is calculated as average of 10 results.

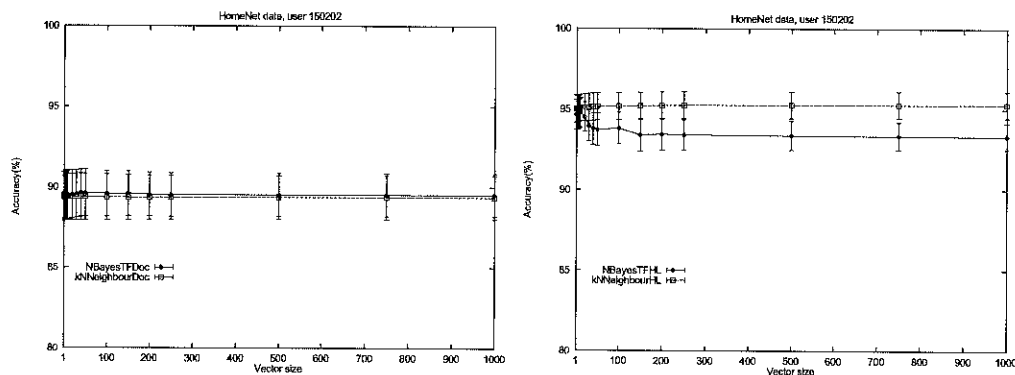


Figure 7: Influence of vector size to classification accuracy of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150202. Notice that classification accuracy scale starts at 80% accuracy. Error bars show standard deviation since accuracy is calculated as average of 10 results.

References

- [1] Apté, C., Damerau, F., Weiss, S.M., Toward Language Independent Automated Learning of Text Categorization Models, *Proc. of the 7th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, Dublin, 1994.

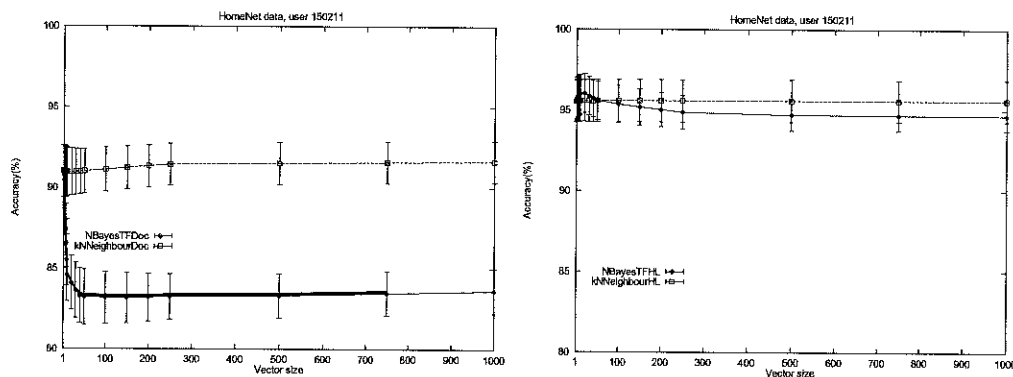


Figure 8: Influence of vector size to classification accuracy of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150211. Notice that classification accuracy scale starts at 80% accuracy. Error bars show standard deviation since accuracy is calculated as average of 10 results.

- [2] Armstrong, R., Freitag, D., Joachims, T., Mitchell, T., WebWatcher: A Learning Apprentice for the World Wide Web, *AAAI 1995 Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, Stanford, March 1995. URL: <http://www.cs.cmu.edu/afs/cs/project/theo-6/web-agent/www/webagents-plus.ps.Z>
- [3] Balabanović, M., Shoham, Y., Learning Information Retrieval Agents: Experiments with Automated Web Browsing, *AAAI 1995 Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, Stanford, March 1995. URL: <http://robotics.stanford.edu/people/marko/papers/lira.ps>
- [4] Bartell, B.T., Cottrell, G.W., Belew, R.K., Latent Semantic Indexing is an Optimal Special Case of Multidimensional Scaling, *Proceedings of the ACM SIG Information Retrieval*, Copenhagen, 1992.
- [5] Berry, M.W., Dumais, S.T., O'Brien, G.W., Using linear algebra for intelligent information retrieval. *SIAM Review*, Vol. 37, No. 4, pp. 573–595, December 1995.
- [6] Burke, R., Hammond, K., Kozlovsky, J., Knowledge-based Information Retrieval for Semi-Structured Text, *Working Notes from AAAI Fall Symposium on AI Applications in Knowledge Navigation and Retrieval*, pp. 19–24, American Association for Artificial Intelligence, 1995.
- [7] Caruana, R., Freitag, D., Greedy Attribute Selection, *Proc. of the 11th International Conference on Machine Learning ICML94*, pp. 28–26, 1994.
- [8] Cohen, W.W., Learning to Classify English Text with ILP Methods, *Workshop on Inductive Logic Programming*, Leuven, September 1995.

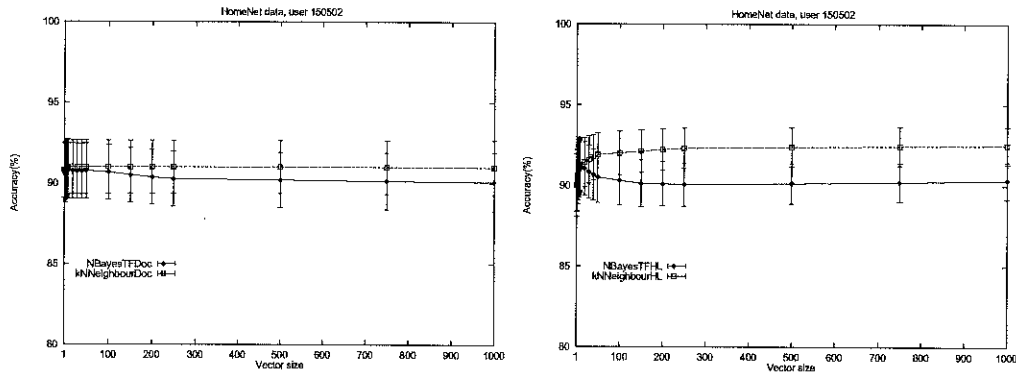


Figure 9: Influence of vector size to classification accuracy of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150502. Notice that classification accuracy scale starts at 80% accuracy. Error bars show standard deviation since accuracy is calculated as average of 10 results.

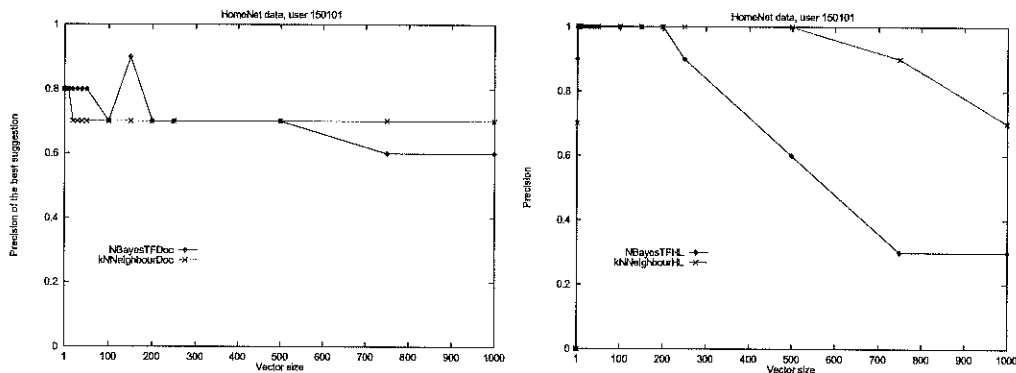


Figure 10: Influence of vector size to precision of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150101.

- [9] Drummond, C., Ionescu, D., Holte, R., A Learning Agent that Assists the Browsing of Software Libraries, *Technical Report TR-95-12*, Computer Science Dept., University of Ottawa, 1995.
- [10] Mc Elligott, M., Sorensen, H., An emergent approach to information filtering, *Abakus. U.C.C. Computer Science Journal*, Vol 1, No. 4, December 1993.
- [11] Etzioni, O., Weld, D., A Softbot-Based Interface to the Internet, *Communications of the ACM* Vol. 37, No. 7, pp.72-79, July 1994.
- [12] Foltz, P. W. and Dumais, S. T., Personalized information delivery: An analysis of information filtering methods. *Communications of the ACM*, 35(12), pp.51-60, 1992.

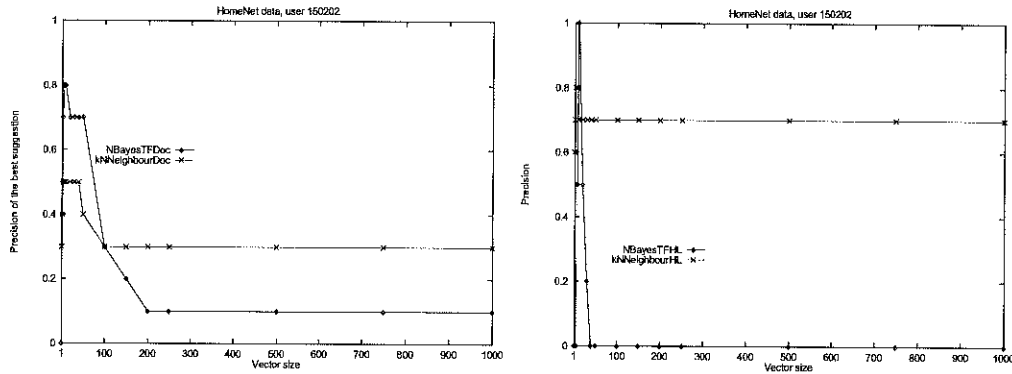


Figure 11: Influence of vector size to precision of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150202.

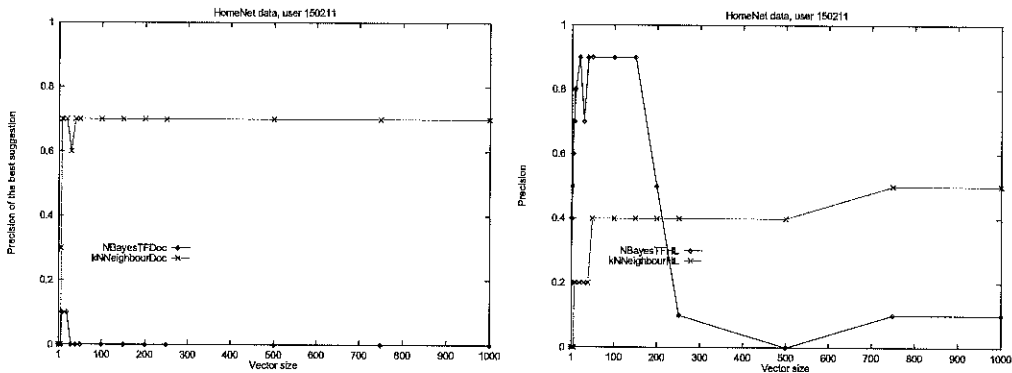


Figure 12: Influence of vector size to precision of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150211.

- [13] Hammond, K., Burke, R., Schmitt, K., A Case-Based Approach to Knowledge Navigation, *AAAI Workshop on Indexing and Reuse in Multimedia Systems*, pp. 46–57, American Association for Artificial Intelligence, 1994.
- [14] HomeNet: a research project at Carnegie Mellon University, Pittsburgh, PA, USA, 1996, URL:<http://homenet.andrew.cmu.edu/progress/>
- [15] Holte, R.C., Drummond, C., A Learning Apprentice For Browsing, *AAAI Spring Symposium on Software Agents*, 1994.
- [16] Joachims, T., Mitchell, T., Freitag, D., Armstrong, R., WebWatcher: Machine Learning and Hypertext, *Fachgruppentreffen Maschinelles Lernen*, Dortmund, August 1995.
- [17] Joachims, T., A Probabilistic Analysis of the Rocchio Algorithm with TFIDF for Text Categorization, *Technical report, CMU-CS-96-118*, School of Computer Science, Carnegie Mellon University, March 1996.

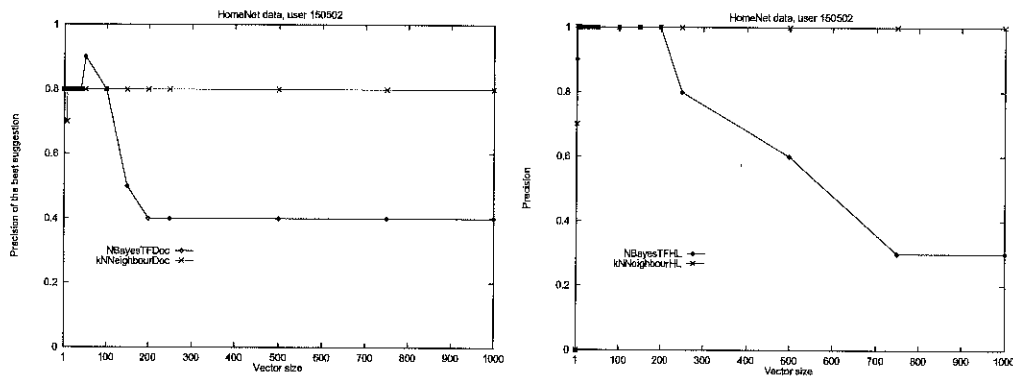


Figure 13: Influence of vector size to precision of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150502.

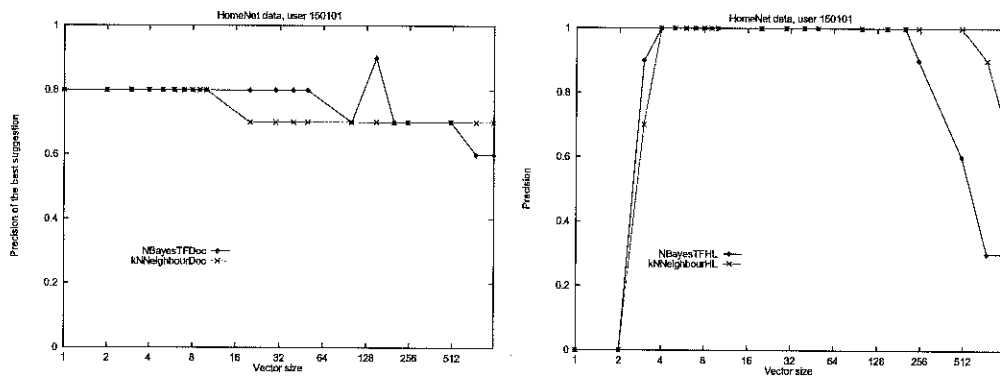


Figure 14: Influence of vector size to precision of model based on documents (left) and hyperlinks (right) for the HomeNet project user with id: 150101. Notice that log-scale on x-axis.

- [18] John, G.H., Kohavi, R., Pfleger, K., Irrelevant Features and the Subset Selection Problem, *Proc. of the 11th International Conference on Machine Learning ICML94*, pp. 121—129, 1994.
- [19] de Kroon, H.C.M, Mitchell, T., Kerckhoffs, E.J.H., Improving Learning Accuracy in Information Filtering, *ICML-96 Workshop Machine learning meets human computer interaction*, 1996. URL: <http://www.ics.forth.gr/~moustaki/ICML96-HCLML/kroon.ps>
- [20] Krulwich, B., Burkey, C., The ContactFinder agent: Answering bulletin board questions with referrals, *Proceedings of the Thirteenth National Conference on Artificial Intelligence AAAI 96*, pp.10—15, 1996.
- [21] Lang, K., News Weeder: Learning to Filter Netnews, *Proc. of the 12th International Conference on Machine Learning ICML95*, 1995.
- [22] Lashkari, Y., The WebHund Personalized Document Filtering System, 1995.

- [23] Lewis, D.,D., Gale, W., A., A Sequential Algorithm for Training Text Classifiers, *Proc. of the 7th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, Dublin, 1994.
- [24] Maes, P., Agents that Reduce Work and Information Overload, *Communications of the ACM* Vol. 37, No. 7, pp.30–40, July 1994.
- [25] Mitchell, T., Caruana, R., Freitag, D., McDermott, J., Zabowski, D., Experience with a Learning Personal Assistant, *Communications of the ACM* Vol. 37, No. 7, pp.81–91, July 1994. URL: <http://www.cs.cmu.edu/afs/cs/user/mitchell/ftp/cacm.ps.Z>
- [26] Pazzani, M., Muramatsu, J., Billsus, D., Syskill & Webert: Identifying interesting web sites, *AAAI Spring Symposium on Machine Learning in Information Access*, Stanford, March 1996 and *Proceedings of the Thirteenth National Conference on Artificial Intelligence AAAI 96*, pp.54–61, 1996.
- [27] van Rijsbergen, C.J., Harper, D.J., Porter, M.F., The selection of good search terms, *Information Processing & Management*, 17, pp.77–91, 1981.
- [28] Shaw Jr, W.M., Term-relevance computations and perfect retrieval performance, *Information Processing & Management*, 31(4), pp.491–498, 1995.
- [29] Rocchio, J., Relevance Feedback in Information Retrieval, in *The SMART Retrieval System: Experiments in Automatic Document Processing*, Chapter 14, pp.313–323, Prentice-Hall Inc., 1971.
- [30] Salton, G., Buckley, C., Term Weighting Approaches in Automatic Text Retrieval, *Technical report, COR-87-881*, Department of Computer Science, Cornell University, November 1987.
- [31] Shardanand, U., Maes, P., Social Information Filtering: Algorithms for Automating “Word of Mouth”, *CHI’95 Mosaic of Creativity*, pp. 210–217, 1995.
- [32] Skalak, D.B., Prototype and Feature Selection by Sampling and Random Mutation Hill Climbing Algorithms, *Proc. of the 11th International Conference on Machine Learning ICML94*, pp. 293–301, 1994.
- [33] Sorensen, H., McElligott, M., PSUN: A Profiling System for Usenet News, *CIKM’95 Intelligent Information Agents Workshop*, Baltimore, December 1995.
- [34] Quinlan, J.R., Constructing Decision Tree in *C4.5: Programs for Machine Learning*, pp.17–26, Morgan Kaufman Publishers, 1993.
- [35] Yang, Y., Expert Network: Effective and Efficient Learning form Human Decisions in Text Categorization and Retrieval, *Proc. of the 7th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, Dublin, 1994.

EXHIBIT D

A Personal Evolvable Advisor for WWW Knowledge-Based Systems

M. Montebello, W.A. Gray, S. Hurley
Computer Science Department
University of Wales, Cardiff.

email: (m.montebello,w.a.gray,s.hurley)@cs.cf.ac.uk

Abstract

The immense size of the distributed WWW knowledge-base and the dramatic rapid increase in the volume of data on the Internet, requires techniques and tools that reduce users' information overload and improve the effectiveness of online information access. Despite the potential benefits of existing indexing, retrieving and searching techniques in assisting users in the browsing process, little has been done to ensure that the information presented is of a high recall and precision standard. In this position paper we present a system that reuses the information generated from search engines together with previously developed systems, and adapts it, by generating user profiles, to better meet the needs and interests of the users by improving recall and precision measures.

Introduction - Background and Motivations

In recent years there has been a well-publicized explosion of information available on the Internet, and a corresponding increase in usage. This is particularly true of the World-Wide Web (WWW) [Berners-Lee et al., 1994] and its associated browsers which allow relative easy access to the information available, and thus make it accessible to a wider audience. The WWW is a major knowledge dissemination system that makes the world's staggering wealth of knowledge and experience, stored on server machines scattered across the internet, accessible to the on-line world.

When people access the web, they are either searching for specific information, or they are simply browsing, looking for something new or interesting (often referred to as *surfing*). The WWW's sheer scale and its exponential growth renders the task of simply finding information, tucked away in some Web site, laborious, tedious, long-winded and time consuming. The fact that a user's time is valuable and that relevant information might not be accessed, imposes serious restrictions on the efficient use of the WWW and the benefits that users can expect from their interaction.

It is well documented that traditional search engines provide services which are far from satisfactory [DeBra and Post, 1994, Spetka, 1994, Srinivasan et al., 1996]. Users are faced with the problem of these search engines being too generalised and not focused enough to their real and specific needs. This triggered further research to develop more sophisticated techniques and agent like systems that make use of the user profile to personalise the service they provide and add value to the information they presented [Pazzani et al., 1996, Green and Edwards, 1996, Mladenec, 1996, Pazzani et al., 1996].

The Personal Evolvable Advisor (PEA), presented in this position paper, is a system we have developed to reuse information generated by search engines and utilise previously developed retrieval systems. Conceptually, the PEA is similar to a meta-search engine, but with the major difference that it employs user profiling to specifically target documents for individual users. In this way duplication and redundancy of information is significantly reduced, while the real needs and interests of the users are fully addressed in a more focussed retrieval.

PEA - Current Implementation

Our goal with PEA is to achieve a high recall and high precision performance score on the information presented to the user. Recall measures how efficient the system is at retrieving the relevant documents from the WWW, while precision measures the relevance of the retrieved set of documents to the user requirements. In order to obtain a high recall execution we make use of the hits returned by a number of traditional search engines together with the output from retrieval systems that have been previously developed. The reason for doing this is twofold. Firstly, we could have developed our own search engine and argued that it utilises the ultimate retrieval techniques and produced results similar to other systems. However, by making use of what other systems generate, we ensure that we obtain all the information that all of them would retrieve at the same time, and not have the problem of developing an ultimate system. Secondly, there are numerous WWW crawlers available, bombarding servers and clogging networks. By using them, we simply use other systems' knowledge-bases, rather than duplicating it, and move up to the next level of the information "food chain" [Selberg and Etzioni, 1995], in this way our recall is as good as can be achieved with any current system. On the other hand in order to add value to the retrieved results and maximise the precision and efficiency with which the system achieves high recall scores, we generate user profiles to predict and suggest the most suitable information for specific users. Through various

interactions the system will be able to optimize the targeting and predicting of what users are interested in, thereby improving the precision factor of the retrieved information.

The processes required by an information retrieval and filtering system include several tasks that PEA decomposes into a number of simpler tasks. Figure 1 shows the major components of the system: the WWW and the external systems at the bottom level, the underlying application software on the next level up, and the GUI at the top.

The WWW is one of the components over which we have no control. It requires no local development, but its heterogeneous, unstructured and uncensored nature causes developers

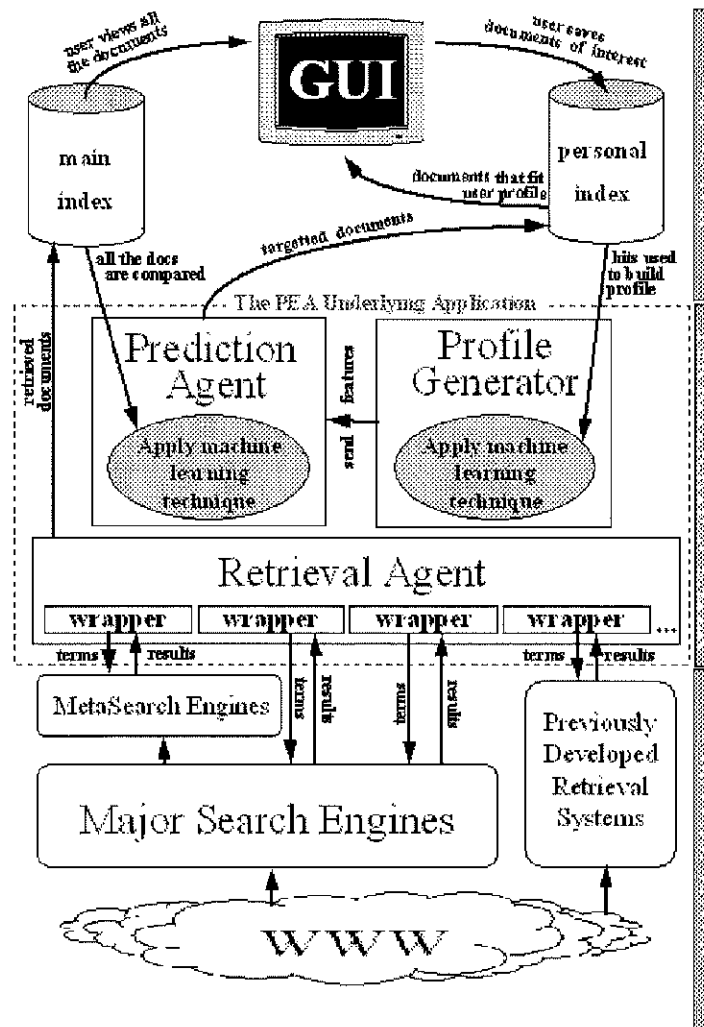


Figure 1: PEA Architecture

to face awkward coding situations in order to be able to cater for all the kinds of data found on the WWW. The WWW can be assumed to be a very large heterogeneous distributed digital information database. In order to optimize the management and exploit the potential of the WWW's vast knowledge-base we require to search and retrieve efficiently and effectively specific information for users.

The external systems utilised include some of the major search engines and also some other retrieving systems that have been developed by other research groups [Eichmann and Wu, 1996, Mladenic, 1996, Selberg and Etzioni, 1995]. They use the WWW as their source of input and we use their output as the input for PEA. All the external systems are considered to be black boxes and action is taken upon the information they output. Wrappers are used to manage the appropriate and proper handshaking between the diverse search engines and the other retrieving systems and the application layer.

Query terms are used to locate documents and retrieve results from the external systems. These results need substantial re-formatting as they usually include completely useless information like advertisements, local links and site specific information.

The underlying application layer has the difficult task of performing all the work required, transparently from the user. It makes use of the information retrieved from the external systems and attempts to improve on the recall/precision metrics mentioned earlier. By using state-of-the-art external systems we attempt to achieve a high recall rate, while using a personalised profile with specific interests for each user, we attempt to also achieve a high precision rate.

The three main components of the PEA underlying application layer (retrieval agent, profile generator and prediction agent - Figure 1) perform the necessary work to satisfy our initial motivations.

The Retrieval Agent

The retrieval agent, is responsible for aggregating all the hits returned by the external systems. It collates the results, by removing duplicates and ensuring integrity, and stores the formatted and pre-ranked results as a single list in a local database, known as the *main index*. The Java programming language was employed to develop this part of the application due to its ease in performing TCP/IP connections to allow retrieval of documents and their processing. This agent interacts with the external systems via appropriate wrappers. Every query term is employed by the wrapper which will command the associated system to locate documents from its local index and return related results. These results are basically a series of document addresses (URLs - Universal Resource Locator) which are listed within an HTML page that the external system returns. A scan through the WWW page will quickly identify the URL links and list them. Some of the links are useless to the user, so the retrieval agent initially removes adverts, duplicates, and site specific links. It then analyses the vetted URLs and accesses the document on-line. This will identify whether the link is still accessible, has moved or been removed completely. If the document is valid, then an initial paragraph from the document is extracted and saved locally in the main database index together with the reference search term, its reference within the index, the URL, and the document title. All these details will be available to the user through the GUI, and also to the prediction agent to identify if the particular document is relevant to a particular user or not.

The Profile Generator

The task that the profile generator sets out to achieve is to analyse each users' personal index and generate a profile. If users have different interests stored in their personal index, then a separate profile is required and generated for each interest. No novel machine learning technique has been developed for the profile generator. It uses specific techniques previously employed by other similar systems [Edwards et al., 1995, Green and Edwards, 1996, Payne and Edwards, 1997]. The difference is that users are able to select which technique they would like to use to generate their profile, and predict other relevant documents in future interactions. This profile generation utilizes the *term frequency/inverse document frequency* machine learning technique [Salton and McGill, 1983], but other machine learning techniques are being implemented. Profile generating systems like MAGI and UNA were relatively easy because specific and fixed fields were provided in the data they were extracting information from. Documents like email and USENET news articles have inherently static field holders embedded in them, e.g. "to", "from", "date", and "subject". These are typical examples of anchored features a developer can rely on when designing the filtering procedures. On the other hand, when considering how to perform the same task on WWW documents (normally HTML), no fixed fields are provided. Even though HTML version 3 introduced the META tag, which allows authors to specify indexing information; it is unreliable as authors can fail to use it. A developer cannot assume that HTML document authors abide by standard conventional fields within documents e.g. "<HTML>", "<TITLE>", and "<BODY>", due to the weak typing nature of HTML. Despite this, there are many systems that filter HTML documents e.g. WebHunter [Lashkari, 1995], LIRA [Balabanovic and Shoham, 1995], Letizia [Lieberman, 1995], WebWatcher [Armstrong et al., 1995] [Joachims et al., 1997], SULLA [Eichmann and Wu, 1996], Personal WebWatcher [Mladenic, 1996], and others described in [Etzioni and Weld, 1994] [Holte and Drummond, 1994] and [Perkowitz and Etzioni, 1995].

We assume that normally, when searching or even browsing, a user bookmarks a page of interest and proceeds with the activity he/she was performing. Taking this activity into perspective, all that is required is to take into consideration what the user bookmarks, and utilise this information to generate the profile. While this method may have problems of over identification, it is more reliable than asking users to assign ratings, as it is less demanding on the user's time. Another problem that many of these HTML filtering systems ignore is that machine learning techniques have a slow learning curve and require a sufficient number of examples before they can make accurate predictions. As a result a profile generator encounters problems when dealing with completely new situations. Generally this is true for all such systems and as [Maes and Kozierok, 1993] rightly argue, the user and the profile agent will gradually build up a trust relationship over time. Issues regarding how many profiles to generate for a user - one specific profile per user, a general profile for a group of users, different profiles for different users or different profiles for the same users - have been tackled differently. Some profile generators develop the 'specific user profile', especially those systems which have been produced to cater for specific items like emails or newsgroups, while others specialise in a 'specific topic profile', like WebFind [Monge and Elkan, 1995], MetaCrawler [Selberg and Etzioni, 1995], PAINT [Oostendorp et al., 1994], and CURRY [Krishnamurthy and Tsangaris, 1996] which recommend documents to users with the same interests or needs. Other systems, like Syskill and Webert [Pazzani et al., 1996], learn a separate profile for each topic of each user. They argue that many users have multiple interests and it will be possible to learn a more accurate profile for each topic separately since the factors that make one topic interesting are unlikely to make another interesting. We take this argument one step further, and argue that what one user finds

interesting in a specific topic, differs from what another user describes as interesting about the same topic. Therefore, different profiles need to be generated for every different interest a user has if the predicted results are to be focused accurately.

The Prediction Agent

The user interest profile generated by the profile generator will be used by the prediction agent in combination with the extracted features from documents in order to predict and suggest new interesting documents to a user. Documents that have been retrieved and stored within the main index by the retrieval agent will have their features extracted and compared to the profile of each individual user generated by the profile generator. This is performed on every item a user has shown interest in, and if any of the documents from the main index happen to fit the user's interests or needs, then they will be eventually suggested to the user the next time the user logs in (Figure 2). Each suggestion, if considered interesting, may be explicitly added to the personal database by the user, or deleted completely. The user might even prefer that he/she is notified, via email, that documents of interest have been located. The machine learning techniques employed to generate the user profile is also applied to extract features from documents. In this way the targetted documents reflect, and are consistent, with the specific user profile generated.

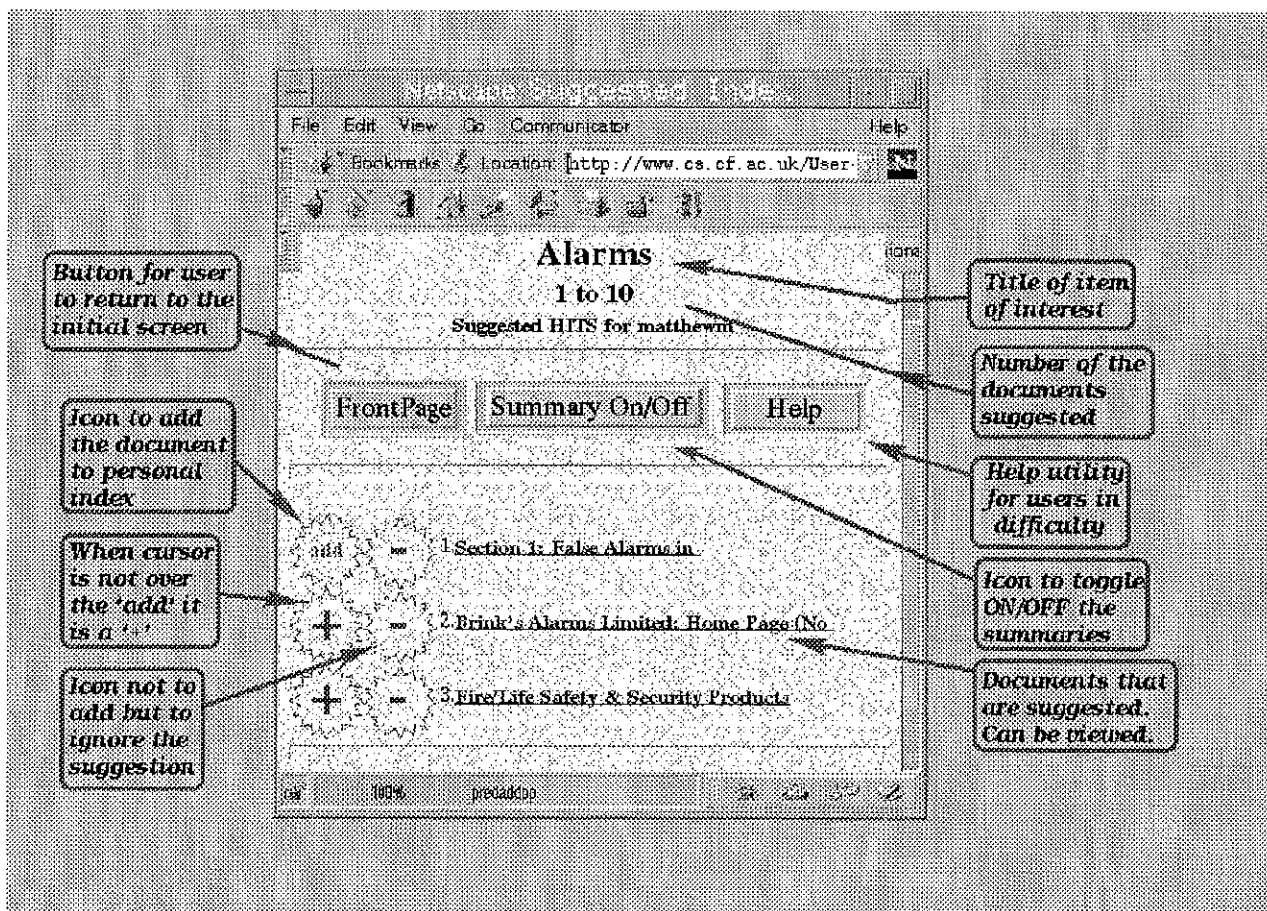


Figure 2: Suggested Documents

Evolvability

One final point to make is about the evolvability of PEA. The choice of external systems coupled with the appropriate wrapper within the retrieval agent, and the machine learning technique employed in the application layer can be selected from a list of systems and techniques incorporated in the system. PEA allows this list to be amended and hence other systems and techniques that might be developed in the future can be easily incorporated. The use of available systems means the system evolves as they evolve, relatively easily. The only amendments to PEA being if a new wrapper is required.

PEA GUI - An Example

PEA requires an administrator to manage the general needs and demands of a specific interest group of users. Search terms tailored to any type of interest group can be initialised by the administrator and furthermore users will be able to suggest any other terms to add to the main search list. Documents relevant to the specific area of interest are retrieved and stored by the underlying application within the main index, and when a user logs-in he/she is able to benefit from the systems' high recall fidelity. Having analysed the documents, individual users can bookmark and highlight specific items as interesting and appealing. These will be saved inside their personal database index. At this stage the underlying application plays another important role in attaining precise targeting of documents to individual users by generating a profile from the personal database index and predicting other documents from within the main index. Users can decide to add the suggested documents to their personal database index or remove them completely. As new and suggested documents are entered in the personal database index the user profile becomes more focused and finely tuned, as a result of which higher precision results will be achieved.

Related Work

Several research systems and commercial off-the-shelf agents have been developed which are similar to our work, but the closest systems are the so called metasearch engines. [Selberg and Etzioni, 1995] fed search terms to six major search engines and made use of the outcome within the MetaCrawler system. A number of front-end metasearch engines have also been introduced on the market, among them Surfbot^[1], WebCompass^[2], WebFerret^[3], and WebSeeker^[4]. These all have very similar capabilities to the MetaCrawler plus additional features such as monitoring specific documents, verifying links and providing relevance ranking. All these related systems are only as reliable as the search engines that they depend upon. This means that their recall score might be very high because a search is done on many of the most popular search engines on the WWW with a single command, potentially retrieving all possible indexed documents. On the other hand, their low precision factor will require the users to check through the documents returned by the metasearchers to identify which ones are of interest. In our system, this task is performed by the profile/prediction components within the underlying application, which combines the optimization of both recall and precision.

Concluding Comments

In this position paper we have presented a system, PEA, that adds value to the information traditional search engines and other metasearch engines generate from the WWW. We argue that by reusing the information output from several retrieving/indexing systems we ensure a high recall score, while generating a specific user profile to predict and target other documents to specific users, we also ensure a high precision score. Users are able to select their own profile generator/prediction agent from a number of alternatives, reflecting different machine learning techniques employed. New techniques can be integrated into this evolvable system by the system administrator, who can also easily maintain the system's resources and update the search terms specific to a user group. In the future we will be investigating the integration of other machine learning techniques that have been developed and employed by other systems. This will help us to evaluate which technique is best suited to cater for the needs of different users. Evaluation of the recall/precision scores is also required to ensure that value is added to the normal services provided by the search engines and the meta-search engines. This will be done by analysing the feedback given from a group of users who are presently making use of the system and who will eventually assess the extent to which the information presented is of high recall/precision quality.

References

Armstrong et al., 1995

Armstrong, R., Freitag, D., Joachims, T., and Mitchell, T. (1995).
WebWatcher: A Learning Apprentice for the World Wide Web.
AAAI Spring Symposium on Information Gathering.

Balabanovic and Shoham, 1995

Balabanovic, M. and Shoham, Y. (1995).
Learning Information Retrieval Agents: Experiments with Automated Web Browsing.
AAAI Spring Symposium on Information Gathering.

Berners-Lee et al., 1994

Berners-Lee, T., Caillan, R., Luotonen, A., Nielsen, H. F., and Secret, A. (1994).
The World-Wide Web.
Communications of the ACM, 37(8):76-82.

DeBra and Post, 1994

DeBra, P. M. E. and Post, R. D. J. (1994).
Searching for arbitrary information in the WWW: the fish-search for mosaic.
In *Proceedings of the 2nd. international world wide web conference*.

Edwards et al., 1995

Edwards, P., Bayer, D., Green, C. L., and Payne, T. R. (1995).
Experience with Learning Agents which Manage Internet-Based Information.
In *Proceedings of ML95 Workshop on Agents that Learn from Other Agents*.

Eichmann and Wu, 1996

Eichmann, D. and Wu, J. (1996).
Sulla - a user agent for the web.
In *5th. International World Wide Web Conference*.

Etzioni and Weld, 1994

Etzioni, O. and Weld, D. S. (1994).
A Softbot-Based Interface to the Internet.
Communications of the ACM, 37(7):72-79.

Green and Edwards, 1996

Green, C. L. and Edwards, P. (1996).
Using Machine Learning to enhance software tools for internet information management.
In Franz, A. and Kitano, H., editors, *AAAI-96, Workshop on Internet-Based Information Systems*, pages 48-55. AAAI Press.

Holte and Drummond, 1994

Holte, R. and Drummond, C. (1994).
A Learning Apprentice for Browsing.
AAAI Spring Symposium on Software Agents.

Joachims et al., 1997

Joachims, T., Mitchell, T., and Freitag, D. (1997).
WebWatcher: A Tour Guide for the World Wide Web.
IJCAI97.

Krishnamurthy and Tsangaris, 1996

Krishnamurthy, B. and Tsangaris, M. (1996).
Curry: A customizable url recommendation repository.
In *5th. International World Wide Web Conference*.

Lashkari, 1995

Lashkari, Y. (1995).
Feature Guided Automated Collaborative Filtering.
Master's thesis, MIT, department of Media Arts and Sciences.

Lieberman, 1995

Lieberman, H. (1995).
Letizia: An Agent that assists Web Browsing.
In *International Joint Conference on Artificial Intelligence*.

Maes and Kozierok, 1993

Maes, P. and Kozierok, R. (1993).
Learning Interface Agents.
In *Proceedings of the 11th National Conference on Artificial Intelligence*, pages 450-465.

Mladenic, 1996

Mladenic, D. (1996).
Personal WebWatcher: Implementation and Design.
Technical Report IJS-DP-7472, J Stefan Institute, Ljubljana, Slovenia.

Monge and Elkan, 1995

Monge, A. E. and Elkan, C. P. (1995).
Integrating external information sources to guide Worldwide Web Information Retrieval.
Technical Report CS96-474, University of California, San Diego.

Oostendorp et al., 1994

Oostendorp, K. A., Punch, W. F., and Wiggins, R. W. (1994).
A tool for individualizing the Web.
In *Proceedings of the 2nd. WWW conference '94: Mosaic and the Web*.

Payne and Edwards, 1997

Payne, T. R. and Edwards, P. (1997).
Interface Agents that Learn: An Investigation of Learning Issues in a Mail Agent Interface.
Applied Artificial Intelligence, 11(1):1-32.

Pazzani et al., 1996

Pazzani, M., Muramatsu, J., and Billsus, D. (1996).
Syskill and Webert: Identifying Interesting Web Sites.
AAAI Conference.

Perkowitz and Etzioni, 1995

Perkowitz, M. and Etzioni, O. (1995).
Category Translation: Learning to understand Information on the Internet.
In *International Joint Conference on Artificial Intelligence*.

Salton and McGill, 1983

Salton, G. and McGill, M. J. (1983).
Introduction to Modern Information Retrieval.
McGraw-Hill.

Selberg and Etzioni, 1995

Selberg, E. and Etzioni, O. (1995).
Multi-service search and comparison using the meta-crawler.
The Web Revolution. Proceedings of the 4th. international world wide web conference.

Spetka, 1994

Spetka, S. (1994).
The TkWWW Robot: Beyond browsing.
In *Proceedings of the 2nd. WWW conference '94: Mosaic and the Web*.

Srinivasan et al., 1996

Srinivasan, P., Ruiz, M. E., and Lam, W. (1996).
An investigation of indexing on the www.
ASIS '96 Annual Meeting of American Society for Information Science., 33:79-83.

Footnotes

...Systems

Research funded by the Radiocommunications Agency, UK.

...Surfbot

<http://www.surflogic.com/>

...WebCompass

<http://www.quarterdeck.com/qdeck/products/webcompass/>

...WebFerret

<http://www.webferret.com/netferret/webferret.htm>

...WebSeeker

<http://www.ffg.com/seeker/>

M Montebello

3/5/1998

EXHIBIT E

**THIS EXHIBIT HAS BEEN
REDACTED IN ITS ENTIRETY**

EXHIBIT F

**THIS EXHIBIT HAS BEEN
REDACTED IN ITS ENTIRETY**