

An Approach and an Eclipse Based Environment for Enhancing the Navigation Structure of Web Sites

Giuseppe Scanniello

*Department of Mathematics and
Computer Science, University of
Basilicata, Italy
giuseppe.scanniello@unibas.it*

Damiano Distante

*Faculty of Economics,
Tel.M.A. University, Italy
damiano.distante@unitelma.it*

Michele Risi

*Department of Mathematics
and Computer Science,
University of Salerno, Italy
mrisi@unisa.it*

Abstract

This paper presents an approach based on information retrieval and clustering techniques for automatically enhancing the navigation structure of a Web site towards improving navigability. The approach increments the set of navigation links provided in each page of the site with a Semantic Navigation Map, i.e. a set of links enabling navigating from a given page to other pages of the site showing similar or related content. The approach uses Latent Semantic Indexing to compute a dissimilarity measure between the pages of the site and a graph-theoretic clustering algorithm to group pages showing similar or related content, according to the calculated dissimilarity measure. AJAX code is finally used to extend each Web page with an associated Semantic Navigation Map. The paper also presents a prototype of a tool developed to support the approach and the results from a case study conducted to assess the validity and feasibility of the approach.

Keywords: *Web site evolution, navigation evolution, reverse engineering, clone detection, clustering, information retrieval, latent semantic indexing, semantic clustering, semantic navigation map.*

1. Introduction

One of the key factors of success of the World Wide Web has been global and simple access to information thanks to an unlimited number links connecting any part of the Web to any other. Everyone using the Web today to publish or access information takes it for granted that any available page will be accessible to anyone who is connected to the Internet [31].

Similarly, the success of a Web site, particularly information intensive Web sites, depends in part on easy and quick access to the information it provides, i.e. on its navigability. The ease of navigation, indeed, is one of the critical factors determining the usability of a Web site, i.e. the capability of the Web site to support the effective, efficient and satisfactory accomplishment of user tasks [40], particularly information gathering tasks.

The importance of navigation in Web sites¹ is also demonstrated by the attention devoted to this aspect by basically all most known Web engineering methods [22], including OOHDM [35], WebML [6], UWE [23], and UWAT+ [14]. All of these methods, indeed, devote a specific design activity and a specific design model to define the navigation structure of a Web site, i.e. the organization of contents into units of consumption and the navigation enabled by links between these units. This design activity is usually named *Navigation (or Hypertext) Design* and the produced model *Navigation Model*. The main modeling concepts a navigation model is usually based on are the concept of *Node* and the concept of *Link*. Nodes are defined as self-contained uniquely identifiable units of information from/to which the user can navigate in a Web site and define links between them. Links are used to connect nodes and to enable navigation between them. Links between nodes are arranged to form particular *access structures* (a.k.a., *navigation patterns*, e.g., *guided tours*, *indexes*, etc.) in order to support specific user information access goals (e.g., quick access to the most important or most recent contents of the site) or specific application requirements (e.g., provide a guided tour to the content of the site on a specific topic) [22].

The adoption of a Web engineering approach to develop a Web site, by organizing the entire development process and by providing developers with the right support in methods, models and tools for design and implementation, helps in producing higher quality Web sites that better satisfy stakeholders' goals and final users' expectations. In particular, when any of these approach is used, it is expected that the resulting Web site will have a navigation structure that, by implementing the designed navigation model, satisfies the users' requirements in terms of contents reachability and navigability.

However, due to time-to-market constraints, lack of proper skills, and/or poor acceptability support [3][18], such Web engineering approaches very often are not used in the industrial practice and few effort is dedicated, in particular, to designing a Web site prior to implementing it. Other times, Web engineering approaches are only used for developing and deploying the first version of the site and then neglected during the rest of the life time of the site when new contents and/or functionalities are introduced and others removed.

¹ Here and elsewhere in the paper, the term "Web site" can be generalized into that of "Web application", as our focus is on their navigation structure.

Both the just described practices may lead, at a certain time, to Web sites suffering of poor usability and incorrect functioning. Regarding navigability, in particular, it may happen that: (i) certain contents become difficult to reach (because of navigation paths too long or too complex) or unreachable (because of missing navigation paths to them); (ii) new published contents may miss links to existing related contents and vice versa; (iii) existing access structures (e.g., guided tours and indexes) may become incomplete, missing to include links towards new related published contents; etc. As a consequence, the user may never get to some of the contents of her/his interest, or s/he may abandon the site because of poor navigability.

It is in these situations, but not only in these, that it may be particularly useful applying the approach proposed in this paper to enhance the navigation structure of a Web site. The approach and the supporting tool we propose enable the automatic extension of the navigation structure of the site by incrementing each page of the site with a sets of links connecting the page with others pages of the site presenting similar or semantically related content. We call this set of links *Semantic Navigation Map*.

As it will be also described later in the paper, another context in which our approach can be profitably applied is that of Web sites built by means of a Content Management System (CMS). Usually, such systems natively support only a category-based access to the contents (often named “articles”) of the built Web sites. Navigation between similar or semantically related contents belonging to different categories has to be supported by means of links expressly defined and kept up to date by the developer or the editor of the site. Our approach, once integrated in a CMS, will support the automatic generation and update of such links.

The approach presented in this paper uses Latent Semantic Indexing (LSI) [11], a well known information retrieval technique, to compute a dissimilarity measure between the pages of a Web site based on the content they present, and a graph-theoretic clustering algorithm [16] to identify cluster (groups) of pages having similar or related content. Links connecting a given page to others pages within the same cluster, i.e., semantic navigation maps, are dynamically injected into each page of the site by means of AJAX code [7]. In order to automate the process of semantic navigation map recovery and injection and to facilitate the adoption of the approach, we have developed a prototype of a supporting tool as an Eclipse plug-in. The approach and the tool have been applied with success in a case study, also presented in this paper, that proved their feasibility and validity.

This paper is an extension of the work presented in [36] and, compared to it, this paper provides the following new contributions:

- A refined version and a more detailed description of the recovery process.
- A detailed description of the tool support, with details on its architecture and implementation technologies.

- An extended version of the case study, which now involves three real-world Web sites.
- An assessment of the validity and feasibility of the approach proved on three real-world Web sites.
- A discussion on the possible applications of the approach.

The rest of the paper is organized as follows: Section 2 discusses a number of works related to Web sites evolution, particularly on navigation restructuring, and to techniques applied in our approach. Section 3 describes the process to recover semantic relations among the contents of a Web site and to enhance its navigation structure with our defined semantic navigation maps. Section 4 briefly presents a prototype of a tool developed to support the approach. Section 5 reports the results from a case study conducted on three real-world Web sites with the purpose of validating the approach and assessing its feasibility and advantages. Finally, Section 6 concludes the paper by providing some final remarks and describing avenues for future work.

2. Related Work

In the last decade, the problem of defining methods and tools for the analysis and evolution of Web sites, and in particular their navigation structure, have been extensively studied [1][4][8][13][15][32][33]. Specific forums have also been created to provide researchers and practitioners with venues for discussing issues and propose solutions on the disciplined evolution of Web-based systems [41]. We synthesize and compare to our some of these works.

A tool for analyzing the navigation structure of a Web site, its evolution during the time, and for identifying possible restructuring interventions to improve navigability has been proposed by Ricca and Tonella in [32]. While supporting the analysis and restructuring of Web sites, the ReWeb tool does not support the automatic application of the proposed changes nor the a restructuring of the navigation structure of the site on a semantic base. In [1], Antoniol *et al.* propose a methodology for reengineering a static Web site. The recovered model is based on the Relationship Management Data Model (RMDM) and the ER+ Model proposed within the Relationship Management Methodology (RMM). The methodology enables also restructuring the navigation structure of a Web site but, differently from our approach, the restructuring is not automatic nor based on similarity between Web pages' content.

Bernardi *et al.* have recently proposed REUWA [2], a process and a supporting tool for the semi-automatic recovery of user-centered conceptual models from existing Web applications, according to the Ubiquitous Web Applications (UWA) design

methodology. Similarly to previous methods, this approach also supports evolution tasks on the navigation structure of a Web site, but at current time changes have to be applied manually and no support for recovering links between pages with similar content is provided.

Other authors have proposed approaches for the model-based evolution of Web applications and in particular of their navigation models. In [17] Garrido *et al.* introduced Web Model Refactorings as behavior preserving transformations for the navigation and presentation models of a Web application aimed at improving its design and external quality. A catalogue of refactorings to improve the quality of a navigation model has also been proposed by Cabot *et al.* in [5]. Both these proposals have not yet a tool support enabling the automatically implementation of the model refactorings into the analyzed application and, in our knowledge, none of the proposed refactoring deals with content similarity. Finally, Lowe and Kong propose in [27] NavOptim, an approach to the evaluation and improvement of Web sites' navigational structures based on the optimization of a navigational effort metric which considers the semantic cohesion between pages and task cases. Semantic cohesion (similarity) between page contents and task cases is calculated using semantic vector spaces but no automatic restructuring of the Web site is supported.

Methods and tools for the analysis and evolution of Web sites rely on several base techniques. These include clustering algorithms, information retrieval and clone detection techniques, refactoring, etc. As it will be described in detail later on in the paper, our approach for automatically enhancing the navigation structure of a Web site uses as base techniques LSI, a well known information retrieval technique, and a graph-theoretic clustering algorithm.

LSI is usually used to abstract concepts from texts (and thus from any document in which a text can be found) and compute a similarity measure between texts, based on their semantics. In our approach we use LSI to compute a similarity measure between the pages of the Web site, and thus identify pages with similar or semantically related contents.

LSI has been widely adopted in different application domains and for different purposes [9][24][25][28][29]. For example, Kuhn *et al.* in [24] describe an approach to group software artifacts using LSI. The approach is language independent and tries to group source code containing similar terms in the comments. Different levels of abstraction to understand the semantics of the code (i.e., methods and classes) are considered. A method based on LSI to compare German literature texts is proposed in [29]. Indeed, the author evaluates whether texts by the same author are alike and can be distinguished from the ones by other authors. Texts by the same author are more alike and tend to form separate clusters. The author also observed that LSI separates prose

and poetry texts in two separate clusters. In our work we use LSI to compute a dissimilarity measure between the pages of a Web site based on the content they show to the final user.

Clustering is a topic analysis technique in the maintenance and evolution of Web applications that is aimed at gathering the entities composing the software system (at different level of abstractions) into meaningful and independent groups. Indeed, a large number of analysis and reverse engineering methods proposed in the literature, including some of those cited above, are based on clustering algorithms [2][9][12][33][34]. For example, different authors have used clustering algorithms to identify Web pages showing similar content and/or having similar HTML structure. Ricca and Tonella [33] enhance the approach proposed by Di Lucca *et al.* in [12] (a pairs of pages are clones if the Levenshtein edit distance [26] between the strings encoding the HTML page structures is zero) by using a hierarchical clustering algorithm to identify clusters of duplicated or similar static pages that could be generalized into a dynamic Web page. Differently from the approach proposed in [12], the distance of cloned pages belonging to the same cluster is not zero. Similarly, in [37] the authors propose a semiautomatic approach based on an agglomerative hierarchical clustering algorithm to identify and align static HTML pages having the same structure and content expressed in different languages. The aligned multilingual pages are then merged into MLHTML pages. De Lucia *et al.* [10] also propose a semiautomatic approach based on the Levenshtein edit distance to compute the similarity of two pages at the structural, content, and scripting code levels. Clones are characterized by a similarity threshold that ranges from 0%, for completely different pages, up to 100%, for identical pages. An approach based on a general process that first compares pages at the structural level (i.e., the Levenshtein edit distance) and then groups them using a competitive clustering algorithm (i.e., Winner Takes All) is proposed by De Lucia *et al.* in [8]. Ricca *et al.* propose in [34] an approach to Web sites understanding based on clustering of client-side HTML pages with similar content. In this work the authors first apply Natural Language Processing (NLP) techniques to associate each page with a set of keywords characterizing it and then they use a hierarchical clustering algorithm to group pages having common keywords and, thus, related content. Keywords are weighted so that more specific and relevant keywords receive a higher score. While we also adopt a clustering algorithm to group pages with similar content, we compute similarity between pages by using LSI instead of NLP and keywords identification and weighting.

In [9] a comparison among clustering algorithms to identify similar pages at the content level is presented. In this work, three variants of the agglomerative clustering algorithm, i.e., a divisive clustering algorithm, k-means, and a competitive clustering algorithm, have been considered. The study reveals that the investigated clustering algorithms generally produce comparable results. To compare pages, an LSI based similarity measure is used.

3. The Process

The process to recover semantic navigation maps from a Web site is schematically represented by the UML activity diagram depicted in Figure 1. In this diagram, rounded rectangles represent process phases and subphases. Overall, the process consists of three main phases: *RecoveringSemanticNavigationMaps*, *InjectingClientSideCode* and *DeployingNewWebSite*. The *RecoveringSemanticNavigationMaps* phase groups pages with similar or related content into clusters. This phase is executed only once and has to be repeated only when the content of the Web site changes because of the addition or removal of pages, or because of the modification of the content of some page. The *InjectingClientSideCode* extends each client side page of the site with the AJAX code enabling the runtime querying (via the interaction with a servlet server side) the database of the identified clusters and the displaying of the semantic navigation map associated to each page. Similarly to the previous phase of the process, this phase is executed only once and has to be repeated only if the content of the Web site changes. The AJAX code injected in each page is instead executed at runtime every time the page is displayed. The *DeployingNewWebSite* phase consists in deploying the enhanced version of the site, i.e. the client side pages extended with the AJAX code and the serverlet used to retrieve data on the recovered clusters of similar pages.

In the following subsections we describe more in detail each of the phases and subphases of the process.

3.1 RecoveringSemanticNavigationMaps

The phase *RecoveringSemanticNavigationMaps* is composed of three subphases: *ComputingDissimilarity*, *GroupingSimilarPages* and *RemovingPages*. *ComputingDissimilarity* extracts the textual content of each page (i.e., the text that is presented to a user) of the analyzed Web site and then computes the dissimilarity between any pairs of pages using a measure based on Latent Semantic Indexing (LSI) [11]. A dissimilarity matrix is produced as output of this activity. This matrix is used by the subphase of *GroupingSimilarPages* to identify groups of pages with similar or related content. The pages included in single clusters, i.e., clusters containing only one page, are discarded by the *RemovingPages* subphase and are not considered in the following iterations of the process. The subphases *GroupingSimilarPages* and *RemovingPages* are repeated until no single clusters remain. Details on these subphases are provided in the following.

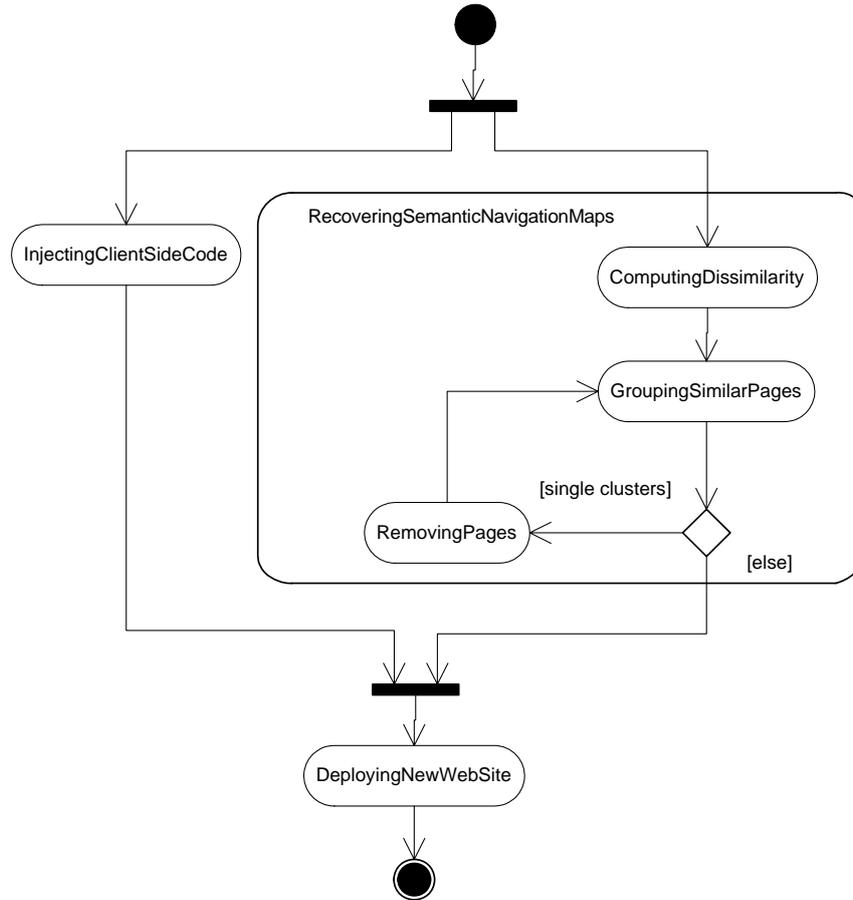


Figure 1. The overall recovery and injection process.

3.1.1 ComputingDissimilarity

This phase is composed of four sequential sub-phases (see Figure 2). *ExtractingPageContent* extracts the textual content of the client-side HTML pages of the given Web site. The extracted content has then to undergo a normalization sub-phase (i.e., *NormalizingContent*) in which non-textual tokens are eliminated (i.e., operators, special symbols, numbers, etc.), terms composed of two or more words are split (e.g., “mail_address” is turned into “mail” and “address”) and terms with a length less than three characters are not considered. A stemming algorithm is also used to reduce inflected (or sometimes derived) terms to their stem. Finally, all the terms contained in a stop word list are removed. The terms within the stop word list should be selected

according to their relevance for the content domain of the considered Web site. Indeed, irrelevant terms should be inserted in that list.

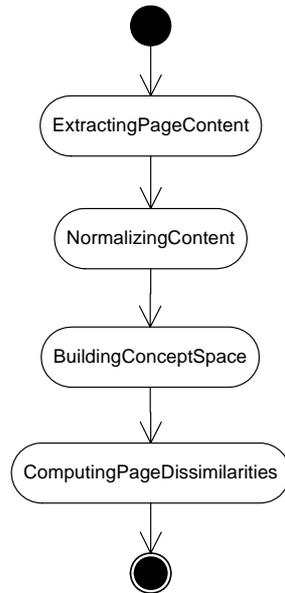


Figure 2. The ComputingDissimilarity’s subphases.

BuildingConceptSpace is in charge of computing the concept space of a Web site on its normalized content adopting LSI. This technique has been originally developed to overcome the synonymy and polysemy problem occurring with the Vector Space Model (VSM) [20]. In fact, LSI explicitly considers dependencies among terms and among documents (corresponding to Web pages in our case), in addition to the associations between terms and documents. This technique assumes that there is a latent structure in word usage that is partially obscured by variability in word choice.

LSI is applied on a term-by-content matrix A , which is built on the normalized content of the considered Web site. In particular, this matrix is $m \times n$, where m is the overall number of different terms appearing in the pages of the site and n is the number of considered pages. An entry $a_{i,j}$ of the term-by-content matrix A represents a measure of the weight of the i -th term in the j -th page. To derive the content latent structure of a Web site we apply on this matrix a Singular Value Decomposition (SVD) [11]. Using this technique the matrix A (having rank r) can be decomposed in the product of three matrices, $T \cdot S \cdot D^T$, where S is an $r \times r$ diagonal matrix of singular values and T and D have orthogonal columns. SVD also provides a simple

strategy for optimal approximate fit using smaller matrices and using only a subset of k concepts corresponding to the largest singular values in S .

The selection of a “good” value of k (i.e., the singular values of the dimensionality reduction of the latent structure) is an open issue. Guidelines to select the suitable value of k have also been proposed in the past, e.g., percentage of number of terms, fixed number of factors, etc. [38]. In our approach, we calculate the number of singular values according to the Guttman-Kaiser criterion [19][21]. This criterion considers the diagonal matrix S of the singular values, which are a kind of eigenvalue, of A in descending order. The value of k is the number of singular values in S more than 1. The rationale for adopting this criterion relies on the fact that it does not require any human intervention, thus fully automating the usage of LSI.

Terms and pages could be graphically represented by vectors in the k space of the underlying concepts of a Web application. In our approach the rows of the reduced matrices of singular vectors are taken as coordinates of points representing the pages in a k dimensional space. To build the dissimilarity matrix of a web application we first compute the cosine between all the pairs of vectors representing the pages in the k dimensional space. The cosine value ranges from -1 (when the two pages have a different semantic) to 1 (when the semantic is the same).

The computation of the dissimilarity matrix requires that the dissimilarities among all the pairs of pages have to be computed. Indeed, the dissimilarity between the pairs of pages is computed in the sub-phase *ComputingPageDissimilarities* normalizing the cosine similarity measure from 0 (when the semantic is the same) to 1 (when they have a different semantic). Hence, given two pages p_1 and p_2 , the dissimilarity between these two pages is defined as:

$$d_{lsi}(p_1, p_2) = \frac{1 - \cos(Vp_1, Vp_2)}{\max_{Vp_i, Vp_j \in W} (1 - \cos(Vp_i, Vp_j))}$$

where Vp_1 and Vp_2 are the vectors corresponding to the pages p_1 and p_2 in the space W of the content of the given Web site. Let us note that this measure cannot be considered a distance as it does not obey the triangle inequality rule. However, this does not influence the possibility of using clustering algorithms as Oudshoff *et al.* show in [30]. Let us note that the rationale for adopting the LSI technique depends on the fact that it has been successfully employed on different applications domains to identify semantic dependences among different entities (e.g., literature texts, software artifacts, source code, web pages, etc.) [9][24][25][28][29].

3.1.2 Grouping Similar Pages

This subphase uses a graph-theoretic clustering algorithm [16] to group pages that are similar at the content level according to the dissimilarity measure described above. Generally, a graph-theoretic clustering algorithm takes as input an undirected graph and then constructs a Minimal Spanning Tree (MST). Clusters are identified pruning the edges of the MST with a weight larger than a given threshold. Nodes within each tree of the obtained forest are included in a cluster.

The graph-theoretic clustering algorithm is used on the strongly connected graph corresponding to the dissimilarity matrix computed in the phase `ComputingDissimilarity`. Nodes are pages and edge weights are the dissimilarity measures between the pairs of pages of the built MST. Clusters are identified pruning the edges with weights less than a given threshold. In our case this threshold is computed as the arithmetic mean of the MST edge weights. We use the arithmetic mean as we aimed at fully automating the clustering process. However, we are conscious that different pruning thresholds may produce better results. Future work is needed to investigate this concern.

In case the graph-theoretic clustering algorithm identifies single clusters, `RemovingPages` is executed. This subphase is in charge of removing from the dataset the pages that the used algorithm includes in the single clusters. The so obtained dataset is then provided again as input to `GroupingSimilarPages` phase. The phases `GroupingSimilarPages` and `RemovingPages` are automatically iterated until no single clusters are identified. This enabled us to improve the overall quality of the clustering process. Furthermore, the overall quality of the clusters (i.e., the ones automatically identified when no single clusters are identified) could be manually improved by modifying the automatically identified clusters. All the pages of a given Web application may be involved in this activity.

Figure 3 shows an example of the clustering results obtained by using the graph-theoretic clustering algorithm on 23 pages of the Web site of the National Gallery of London (i.e., one of the Web site used as case study). In particular, this figure shows the built MST and the clusters identified using 0.33 as threshold. The clustering algorithm identified 9 different groups of pages with similar or related content. Among the identified clusters, six were single clusters. Note that neither the names nor the title of the pages are reported in the figure for readability reasons. Edge weights less than the threshold values are not shown as well.

3.1.3 Removing Pages

As mentioned above, the subphase RemovingPages is executed in case the graph-theoretic clustering algorithm identifies one or more single clusters. This phase is in charge of removing, from the dissimilarity matrix, the rows and the columns corresponding to the pages that have been placed within single clusters. The obtained matrix is then provided as input to GroupingSimilarPages. RemovingPages is executed until no single clusters are identified by the clustering algorithm. The motivation for removing the pages of the single clusters lies in the fact that they should have a low level of semantic similarity with the remaining pages.

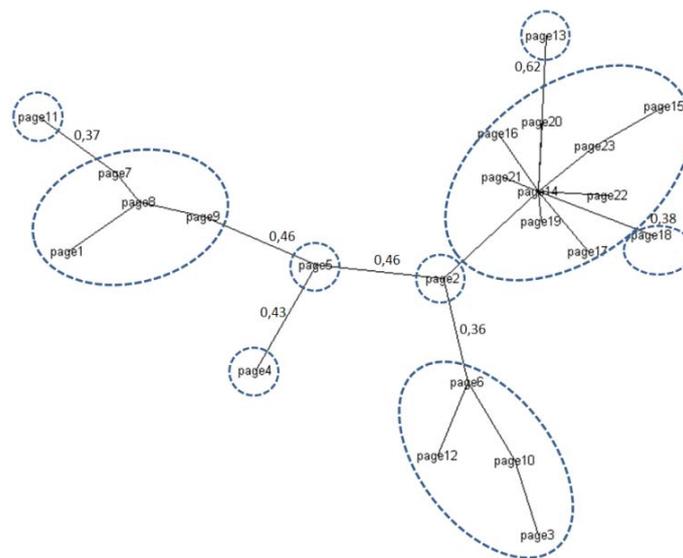


Figure 3. Example of clusters of pages obtained for the NationalGallery.org.uk case study.

3.2 Injecting Code into Client Pages

This phase actually enhances the navigation structure of the considered Web site by introducing a semantic navigation map within each of its pages. Indeed, each page is extended to dynamically include a set of links towards the pages that the clustering process has placed in the same cluster. To this end, we use AJAX and DOM technologies to respectively query the server that maintains the results of the clustering process and to display the semantic navigation maps within the pages of the Web site.

For the first objective we use a Javascript function that retrieves the list of pages within the cluster previously computed of a given page. Figure 4 shows the Javascript code used to get this list. In particular, the Javascript function establishes a connection with the servlet *retrieveTrackingMap* to get the list of pages to visualize within the semantic navigation map.

For the second purpose, we use a different Javascript function (see Figure 5) to interpret the server answer and to properly display the map of links in the considered Web page by injecting HTML code (see Figure 6) into the DOM of the page. To inject the HTML code we add the DIV tag at the end of each page. This code is different for each page of the Web site as it includes the name of the page for which it has been generated. For example, the code shown in Figure 6 regards the page *michelangelo.html*.

```
function getCluster(divBlock,pageCode){
  mapList = divBlock;
  var req = newXMLHttpRequest();

  req.onreadystatechange = getReadyStateHandler(req, displayMap);
  req.open("POST", "<%=request.getContextPath()%>/servlet/retrieveTrackingMap", true);
  req.setRequestHeader("Content-Type", "application/x-www-form-urlencoded");
  req.send("code=" + pageCode);
}
```

Figure 4. The Javascript code used to get data on a page cluster.

```
function displayMap(rdftXML){
  var rdfs = rdftXML.getElementsByTagName("rdf")[0];
  var items = rdfs.getElementsByTagName("item");
  mapList.innerHTML = '';
  var code = '';
  for (var i = 0; i < items.length; i++){
    code +=
      items[i].getElementsByTagName("rank")[0].
      firstChild.nodeValue + "<a href=\" +
      items[i].getElementsByTagName("link")[0].
      firstChild.nodeValue + \">\" + "<b>\" +
      items[i].getElementsByTagName("title")[0].
      firstChild.nodeValue + "</b>\" + "</a><br/>\" +
      items[i].getElementsByTagName("desc")[0].
      firstChild.nodeValue;
  }
  code = createMultiTab(code, 8);
  mapList.innerHTML = code;
  activateMultiTab(mapList);
}
```

Figure 5. The Javascript code used to visualize a semantic navigation map.

At run-time a Javascript code is used to overlap the semantic navigation map on the original web page. This has the effect of displaying the map on the right hand side of the page. However, the end user can place the map wherever s/he wants (see Figure 11 and Figure 12) using drag and drop functionalities. The semantic navigation map preserves its position also when the user

scrolls the page. Furthermore, if necessary, the map can be folded. Note also that the semantic navigation map will provide a multi-tab visualization in case it contains a number of links larger than 8 (see Figure 12).

All the functions required for managing the communication with the server, parsing messages from it, and visualizing the semantic navigation maps, have been collected in a Javascript library (i.e., `localTrackingMap.js`), thus reusing them in each page of the evolved Web site.

```
<!-- Begin Semantic Navigation Map -->
<div class="LocalTrackingMap">
  <div id="localmap"></div>
  <script src="script/localTrackingMap.js" type="text/javascript"></script>
  <script>getCluster("localmap","collection\artist\michelangelo.html");</script>
</div>
<!-- End Semantic Navigation Map -->
```

Figure 6. The HTML code injected in the Web pages.

Note that the HTML pages of a given Web site are modified once for all. This is possible because the clusters are independently detected (see Figure 1) and are dynamically obtained querying the server. Furthermore, the identification of pages with similar or related content should be performed when the content of existing pages is modified or new pages are deployed in the Web site. For consistency reasons, the identification of similar pages should also be repeated when pages are removed.

3.3 Deploying the Enhanced Web Site

In the *DeployingNewWebSite* phase the enhanced pages and the data on the recovered clusters are deployed on the server. The deployed pages require a server component that retrieves data on the semantic navigation map recovered and associated to each client side page. To enable the communication between each enhanced client side page and the server where the corresponding cluster is stored, we have developed a servlet (i.e., `retrieveTrackingMap`). The software engineer has to deploy this servlet only once. Finally, when the enhanced pages and the servlet have been deployed, the Web site is enabled to be accessed by the users.

In this phase, the software engineer should also perform testing to find possible faults of the new version of the Web site. Currently, our tool prototype does not provide any specific support for testing. Future work will be devoted to support the software engineer to execute test cases on the Web pages enhanced with the semantic navigation maps. To this end a Web browser extension could be developed to record, edit, and debug test cases.

4 The Tool Prototype

To support all the phases of the proposed recovery and extension process, we have implemented a prototype of a supporting tool as an Eclipse plug-in. In the following subsections we first present the logical architecture of the tool prototype and then we describe how it supports the software engineer in the different phases of the process.

4.1 Tool Prototype Architecture

Figure 7 shows the layered architecture (i.e., the logical architecture) of the system prototype through a UML Package Diagram. The *Data* component contains the pages of the Web site and all the intermediate data produced by the process phases. The *Computing Dissimilarity* component uses the *HTML Parser* component to extract the content of the HTML client-side pages. HTML Parser integrates an open source HTML parser written in Java (HTMLParser ver. 1.6), available under the GPL license at sourceforge.net².

The content extracted from each page of the site is stored in the local file system (i.e., the Eclipse workspace) to avoid repeating this operation for the same page more than once. The Computing Dissimilarity module uses the *LSI Engine* component to compute the dissimilarity matrix of the page content. To this end this component integrates an R³ implementation of the LSI information retrieval technique. This implementation is available under GPL license from www.cran.r-project.org. The rationale for using R lies in the fact that it is open source and a huge and very active community works to improve it.

The dissimilarity matrix produced by the Computing Dissimilarity component is stored in the Eclipse workspace. This choice is motivated by the fact that the computation of this matrix is expensive and it is performed only in case new pages are added or remove from the original Web site (see Section 3.2). For example, the computation of the dissimilarity matrix for the National Gallery Web site, which included 2017 pages, took about three hours using a laptop equipped with a 1.5 GHz Intel Centrino with 1.5 GB of RAM, a 60GB Hard Disk and Windows XP Professional SP 3 as operating system.

²HTMLParser ver. 1.6 can be downloaded from at <http://sourceforge.net/projects/htmlparser>.

³R is a free software environment for statistical computing and graphics. It is worth noting that there is a very proficient and active community that evolves the environment and the majority of the available libraries are available under GPL license.

The dissimilarity matrix is successively provided as input to the *Grouping Pages* component. This component is in charge of detecting pages that are similar. To this end an R implementation of the graph-theoretic clustering algorithm (available under GPL license from cran.r-project.org) has been integrated. Grouping Pages also removes from the dissimilarity matrix rows and columns corresponding to the pages within single clusters. Note also that the R implementations of the LSI engine and the clustering algorithm are all integrated within the system prototype using Rserve (from cran.r-project.org), a TCP/IP server allowing programs to use R facilities. Despite the availability of simpler methods to integrate R software components with Java code, we decide to use Rserve in order to eventually distribute the computation on different nodes on the Web. This was due to the fact that the time needed to execute the tool on large sized Web sites could be considerable.

Modifying Web Pages injects the needed HTML and Javascript code into the pages of the Web application to enable the communication with the server and to display the semantic navigation maps.

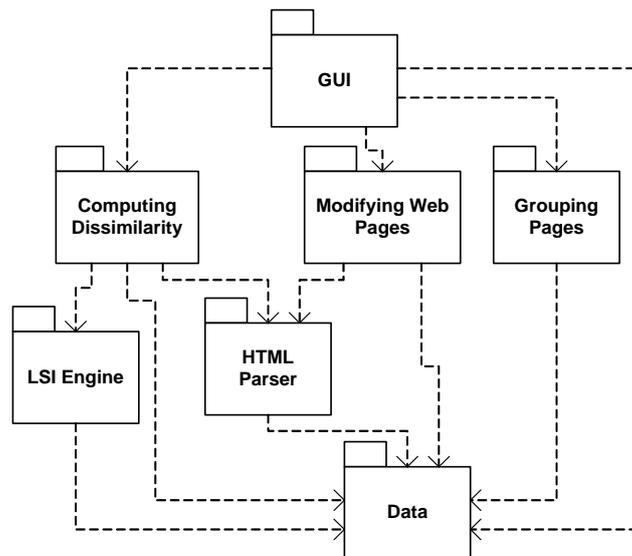


Figure 7. The architecture of the tool prototype.

The point where the HTML and Javascript code has to be injected is identified using the HTML Parser component. The *GUI* component enables the software engineer to select the pages of the Web site and to display the groups of similar pages detected by the tool (see Figure 9). Pages can be manually added or removed from the identified clusters, thus improving the overall quality of the identified groups of similar pages.

4.2 Using the Eclipse Plug-in

The developed Eclipse plug-in fully supports all the phases of the process depicted in Figure 1. As a first operation, the plug-in requires that the software engineer creates an Eclipse project. To this end, the plug-in proposes a wizard, which allows specifying the name of the project and its workspace. Once the project has been created, the analyst has to select the Web site to enhance. This is possible by right-clicking the project within the *Package Explorer* view and choosing *Select Source Directory* within the menu *semantic navigation map* (see Figure 8). The plug-in will provide a window for choosing the path of the folder where the pages of the original Web site are stored. Note that these pages have to be stored on the file system of the personal computer where the plug-in is installed. The plug-in could be extended to work on Web pages maintained on remote servers. This is a further direction to extend the proposed supporting tool.

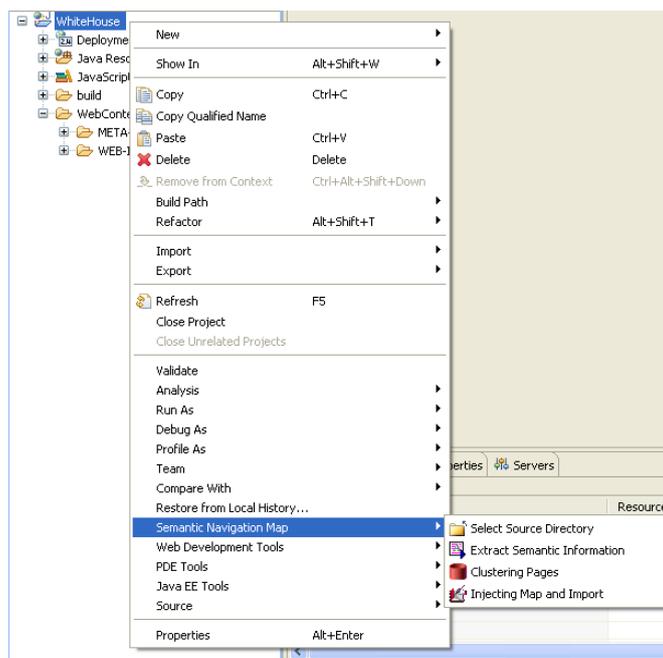


Figure 8. Semantic Navigation Map Menu

The clustering process starts by right-clicking the project and choosing *Clustering Pages* within the menu *semantic navigation map* (see Figure 8). The plug-in will produce a file (i.e., dbcluster.txt) containing all the automatically identified clusters. Indeed, it also contains all the pairs of pages composing each cluster and their similarity level (i.e., the cosine between the vectors of the pages in the content space). To improve the overall quality of the clustering process, the automatically

identified clusters can be manually refined within the plug-in. If needed, the analyst can refine a cluster, by modifying the set of included pages and their similarity levels. Figure 9 shows the report of the clustering results visualized within the plug-in. Note that in the current version of the plug-in the analyst can only edit the report containing the identified clusters by simply using the plug-in text editor. The plug-in does not perform any check on the correctness of the modifications made on the report.

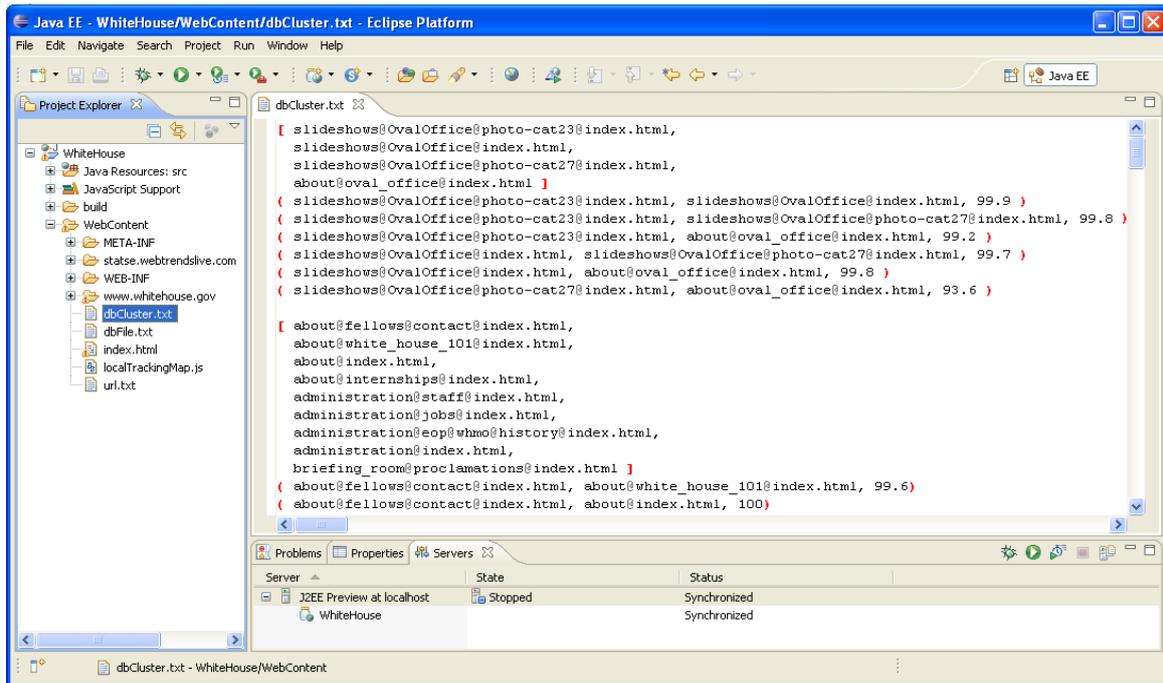


Figure 9. Improving Groups of Similar Pages

The software engineer uses then the Eclipse plug-in to extend the pages of the Web site by automatically introducing the semantic navigation maps (see Figure 8). Successively, the new pages and the Javascript library will be automatically added to the project workspace. Finally, the plug-in enables the deployment of the enhanced Web site on a suitable Web server.

5. Case Study

The approach and the tool prototype have been assessed on three real-world Web sites: the National Gallery of London

(NGL)⁴, Play Shakespeare (PS)⁵, and the White House (WH)⁶. In the following subsections we present and discuss the results obtained from this case study.

5.1 Context

Some pages of NGL presented a navigation menu (see for example Figure 11) to promote a quick access to the Web site content. As the careful reader may object that in the case in point the recovered semantic navigation maps may result in duplicate navigation structures (the navigation menu) and may have biased the obtained results, we have also investigated the effectiveness of the approach on the Web sites PS and WH. The rationale for choosing PS and WH relies on the fact that their pages lacked of a navigational menu that could positively affect the results of the analysis.

To get the pages of the selected Web sites we used a freeware dumper (i.e., HTTrack Website Copier, available at www.httrack.com). The motivation for performing the dump is that the source HTML files were not available for the download. Regarding NGL the dump was executed on June 9th, 2008. In particular, HTTrack Website Copier downloaded 6573 HTML pages using the index page as starting page and following the hyperlinks connecting the pages stored until the sixth level of depth was reached in the folder hierarchy of the Web site. The mirror of NGL has been successively analyzed to prune duplicated HTML pages created by the dumper and documents currently not considered by the approach (e.g., PDF and Word files) and multimedia objects (e.g., JPG images and flash animations). Regarding the HTML pages, we limited the analysis to the ones presenting information and content on the museum entire permanent collection and long term loans. The pages of the works of art shown during the years in the gallery have also been considered (i.e., the pages within the Exhibition section). The total number of HTML pages we have selected and considered in the presented case study is 2017.

The dump of the PS Web site was executed on January 16th, 2009. The dumper downloaded 5570 pages starting from the index page and following all the hyperlinks connecting the pages until the fifth level of depth. Similarly to the NGL Web site, downloaded pages have been analyzed to prune duplicated and meaningless pages, irrelevant documents, and multimedia objects. As a consequence of the pruning, the number of HTML pages was 1770.

Regarding the WH Web site, the dump was performed on January 22nd, 2009. The number of downloaded pages was 424 following all the hyperlinks connecting the index page until the fifth depth level. Also, in this case a pruning phase has been

⁴ www.nationalgallery.org.uk

⁵ www.playshakespeare.com

⁶ www.whitehouse.gov

performed to remove duplicated and meaningless pages, irrelevant documents, and multimedia objects that we do not treat in the approach. Once the pruning was performed the number of remaining pages 313.

5.2 Results

For each analyzed Web site, the selected pages were then provided as input to the plug-in for the following clustering and semantic maps recovery and injection. In the case on NGA, the 2017 pages were grouped into 243 different clusters. The largest cluster contained 32 pages and included pages from the collection *Scientific Instruments and Inventions from the Past* of the gallery. On the other hand, the mean number of pages within the clusters was 2.6. Regarding the PS Web site, 141 were the identified clusters, with an average of pages per cluster equals to 8 and the large cluster counting 436 pages. For the WH Web site, the 313 analyzed pages were grouped into 17 clusters having a mean of 5.5 pages and the large cluster composed of 25 pages. These and other descriptive statistics on the case study are summarized in Table 1. The cardinality distribution of the clusters containing at least two pages for each analyzed Web site is graphically summarized by the boxplots shown in Figure 10. Note that the median for each Web site is either 2 or 3. Overall, even clusters containing 2-3 pages are useful and meaningful, as each of them identifies a set of 2-3 pages having similar and correlated content.

Table 1. Descriptive statistics of the analyzed Web sites.

	NGL	PS	WH
Number of analyzed pages	2017	1770	313
Number of identified clusters	243	141	17
Number of iterations	5	4	5
Number of pages within single clusters	1387	639	219
Percentage of pages placed in single clusters	68%	36%	70%
Number of pages within clusters containing at least two pages	630	1131	94
Number of pages within the largest cluster	32	436	25
Mean number of pages within the clusters	2.6	8.0	5.5
Number of characters within the analyzed pages	3.460K	13.716K	1.853K
k-value	1719	1503	273

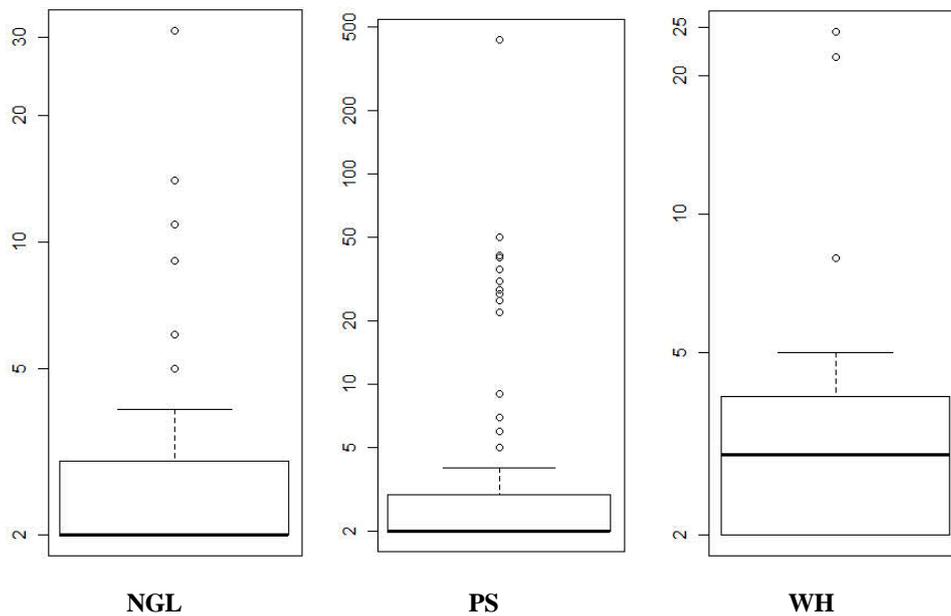


Figure 10. Cardinality distribution of the clusters for the analyzed Web sites.

The Eclipse plug-in was finally used to inject the semantic navigation map code into the pages of all the considered Web sites. Examples of the enhanced version of the pages obtained applying the approach on the Web sites NGL and WH are shown in Figure 11 and Figure 12, respectively. It is worth noting that the new Web sites have been deployed on a local machine of one of the authors. This was due to the fact that we assessed the approach and the system prototype on the local copy of the selected Web sites. Let us also note that the pages show in Figure 11 and Figure 12 differ from their original versions in the presence of the semantic navigation maps (see the right hand sides of the pages). Each semantic navigation map presents a set of links towards pages that have been found similar in content. Each link shows:

- (i) the title of the target page;
- (ii) the percentage of similarity with the current page obtained using LSI;
- (iii) a description of the target page obtained from data in its description meta tag, if available.

THE NATIONAL GALLERY
Trafalgar Square London

Collection

Search: Go Site Map

HOME
COLLECTION
Beginner's Guides
Collection at a Glance
Artists at a Glance
Collection Explorer
Full Collection Index
Collection Features
Collection News
Study & Care
National Inventory
Picture Library
EXHIBITIONS
WHAT'S ON
PLAN YOUR VISIT
ABOUT THE GALLERY
EDUCATION

ON-LINE SHOP SITE
JOBS
SUPPORT THE GALLERY
CONTACT US



The Virgin Mary with the Apostles and Other Saints
1423-4
ANGELICO, Fra
Died: 1455
NG663.2. Bought, 1860.

The figures have been identified as the Virgin Mary (top right) with apostles and the evangelists.

This, along with four other panels showing respectively, 'Christ Glorified in the Court of Heaven', 'The Forerunners of Christ with Saints and Martyrs', 'The Dominican Blessed' and another panel of 'The Dominican Blessed', formed the predella, or lower section, of the high altarpiece of San Domenico at Fiesole, near Florence. This was the church of Fra Angelico's own Dominican friary. The predella shows the most elaborate depiction of the Court of Heaven in the Collection. Christ stands in the centre surrounded by angels, saints and martyrs.

The church of San Domenico was dedicated in 1435, and Fra Angelico's picture was probably in place on the high altar by that time. The main panel was modified by Lorenzo di Credi around 1501. This and the painted plaster are still in the church.

Egg tempera on wood
32 x 64 cm.

Online Shop
Buy a Print of this Painting
Books & Catalogues

See Also...
Other works by this artist...
Christ Glorified in the Court of Heaven
The Dominican Blessed
The Dominican Blessed
The Forerunners of Christ with Saints and Martyrs

Show descriptions Scrollable Draggable Expanded
Suggested related pages (5)

- 97.9% **The Dominican Blessed**
- 97.9% **The Dominican Blessed**
- 97.6% **The Forerunners of Christ with Saints and Martyrs**
- 96.6% **Christ Glorified in the Court of Heaven**
- 89.8% **Selected Altarpieces 1260-1450**

Figure 11. The enhanced version of a page from the National Gallery of London Web site.

the WHITE HOUSE PRESIDENT BARACK OBAMA

the BRIEFING ROOM ISSUES the ADMINISTRATION ABOUT the WHITE HOUSE our GOVERNMENT CONTACT us

CHANGE HAS COME TO AMERICA

THE INAUGURATION of PRESIDENT BARACK OBAMA

Watch video from the Inaugural Ceremonies and read President Obama's Inaugural Address.

WATCH NOW

MORE FEATURES

1 2 3 4



THE BLOG MORE FROM THE BLOG

WED, JANUARY 21, 1:27 PM EST
President Barack Obama's Inaugural Address
Yesterday, President Obama delivered his Inaugural Address, calling for a "new era of responsibility."
Watch the video.
READ THIS POST

TUE, JANUARY 20, 2:15 PM EST
A National Day of Renewal and Reconciliation
President Barack Obama's first proclamation.
READ THIS POST

TUE, JANUARY 20, 12:01 PM EST
Change has come to WhiteHouse.gov
The first post on WhiteHouse.gov.
READ THIS POST

SEARCH THE SITE

AGENDA

Economy
Read the President's economic agenda.

Energy & the Environment
Read the President's agenda on energy & the environment.

More Issues
Read the President's entire agenda.

WHISTLESTOP TRAIN TOUR

WATCH NOW

Show descriptions Scrollable Draggable Expanded
Suggested related pages (24)

1-8 9-16 17-24

- 94.2% The Briefing Room
- 92.2% Blog
- 87.4% The Administration
- 87.2% White House Jobs
- 86.9% Weekly Video Address
- 86.2% The White House - Blog Post - The Whistle Stop Tour
- 85.7% The Agenda
- 85.6% Proclamations

Figure 12. The enhanced version of a page from the White House Web site.

5.3 Discussion

By analyzing the pages of the Collection section of the NGL Web site (i.e., the pages presenting information on the permanent collection and long term loans of the museum) we noted the presence of a navigation menu on the bottom right hand side (see Figure 11). In particular, given an artist and his/her work of art, the menu enables users to directly access the pages presenting other works of art of the same artist exhibited at the National Gallery of London. Additional links are also proposed to navigate and to access related pages. Since NGL presented a menu, we compared it with the recovered semantic navigation maps to assess the effectiveness of the approach. This comparison revealed that most of the links presented by the navigation menu were proposed by analogous links in the recovered semantic navigation map. As an example, if we compare the navigation menu and the navigation map of the page shown in Figure 11 we can observe that they show basically the same links. The only difference we can notice is the link *Selected Altarpieces 1260-1450* that is present in the semantic navigation map and not in the navigation menu. This link allows accessing a page showing the works of art from different artists on the theme of altarpieces. On the other hand, we also noted that some pages of the Collection section presented a navigation menu with a number of links larger than the ones within semantic navigation maps. However, the maps of the majority of these pages present the same links as the corresponding navigational menu. This indicates that the links of the identified semantic navigation maps are correct.

Even on the PS Web site the plug-in produced interesting results. We observed that the majority of the links included in the semantic navigation maps connected pages with similar or semantically related content. For example, on PS, the larger identified cluster mainly contained all the pages presenting information regarding the characters of the Shakespeare's plays. Indeed, 327 were the pages within this cluster presenting information on the play characters. Overall, we can observe that also clusters containing 2 or 3 pages are useful and meaningful, as each of them identifies a set of 2 or 3 pages, out of the totality of analyzed pages, having similar or correlated content.

Due to the size of the considered Web sites we cannot analyze the completeness of the identified maps. However, some considerations can be made according to the descriptive statistics presented in Table 1. In particular, on PS and WH the average number of pages within the identified clusters seems correct. Also, the number of pages within clusters containing at least two pages enables us to believe that the obtained results could be appropriate. A further observation is motivated by the number of single clusters identified by the Eclipse plug-in. This could bias the usefulness of the approach as single clusters do not contribute to the enhancement of the navigational structure of a Web site. Hence, the larger is the percentage of pages placed in

clusters containing at least two pages the more is the usefulness of the approach. In fact, we observed that on the Web sites NGL and WH the percentage values of the single clusters with respect the total number of analyzed pages (see Table 1) were 68% and 70%, respectively. On the other hand, a smaller percentage value (i.e., 36%) was achieved on PS. Such a difference was due to the fact that PS presents pages richer in content compared to the other two Web sites.

Overall, the results obtained on the Web sites selected as case study suggest three considerations. The first concerns the subjective satisfaction of the users. This issue will be addressed in the future conducting special designed investigations, e.g., survey questionnaires and interviews [22][39]. The latter two considerations regard correctness (links included in the navigation maps actually propose pages with content similar or related to the one showed by the considered page) and completeness (the list of proposed links includes most of the pages with actually related content) that will be discussed in the following subsection.

5.3.1 Assessing correctness and completeness

The correctness and the completeness of the approach have been assessed on the static pages of SRA (Student Route Analysis) a dynamic Web site implemented by one of the author. SRA was developed to provide the students in Politics Science at the University of Salerno with statistics and information on their academic carrier. SRA was composed of 380 files distributed in 45 folders according to a meaningful classification. Overall, the SRA web application was composed of 24 html pages and 45 dynamic pages.

To quantitatively assess the produced results we used two well known metrics, namely precision and recall. The precision and recall have been used to assess the correctness and the completeness of the automatic identified clusters, respectively. In our case the precision is the number of actual pairs of similar pages identified by the tool over the total number of identified pairs, while the recall is the ratio between the number of actual pairs of similar pages identified by the tool over the total number of actual pairs of similar pages.

The obtained precision value was 0.89, while the recall was 0.57. This indicates that almost all the links included in each semantic navigation map are correct. Moreover, about 2/3 of the actual links between pairs of pages have been correctly detected. Although the achieved results are encouraging a further investigation is needed to further assess the effectiveness of our approach in terms of correctness and completeness.

6. Conclusion and Future Work

Our research is meant to automatically recover semantic relations between the contents of a Web site and build semantic navigation maps, accordingly. To this aim, we have defined a process that first computes the dissimilarity between Web pages using a measure based on Latent Semantic Indexing and then uses the computed dissimilarity to group pages showing similar or related content, by means of a graph-theoretic clustering algorithm. Finally, the navigation structure of the Web site is enhanced by adding links between pages within the same cluster. We defined the set of links added to a page to connect it to other pages of the same cluster as Semantic Navigation Map. To automate the application of the approach, we have developed a supporting tool as an Eclipse plug-in. Both the approach and the tool have been validated in a case study involving real-world Web sites.

When we have started this work, we considered the client-side HTML pages of the site as the elementary granules of information (nodes) to analyze, index and link, and assumed that the content (text) included in a page and showed to the final user is uniquely identified by the page URL. Therefore, the proposed process and the supporting tool are applicable to Web sites in which the content associated to a page does not depend on the user logged on, nor on other context variables, though it may change in time⁷. Our approach is instead applicable no matter of the technologies used server-side to produce the front-end of the application, provided that this is represented by HTML pages.

Semantic navigation maps may be particularly useful when the navigation structure of the site is found to be not properly designed or when it has degraded during the Web site life-time. However, even properly designed Web sites may benefit from the use of the semantic navigation maps. In fact, they represent an additional navigation structure that is complementary to those implementing the navigation model obtained during the design phase. While the navigation model of a Web site is usually designed to satisfy specific navigation requirements (e.g., provide access to subsets of contents grouped by category or having some characteristic of interest for the user, etc.), our semantic navigation maps are intended to make explicit the latent semantic relations between the contents of the site and keep this relation up-to-date when the content of the site evolves.

Similarly to Web site search engines, such as Google Site Search (available at www.google.com/coop/cse) or FreeFind (available at www.freefind.com), semantic navigation maps represent a navigation structure built by analyzing and indexing the contents of the Web site based on their meaning. Differently from them, the index provided by a semantic navigation map is not

⁷ A similar limitation affects also search engines, which are able to index only the freely accessible (portion of) Web sites.

the result of a search query executed by the user at run-time, but a navigation structure provided by default by the site to its users. More importantly, the navigation map is obtained considering as “input parameters” of the search query not a particular keyword or set of keywords, but the full text of the page the user is visiting. The same consideration applies if we compare the semantic navigation maps recovered by our approach to the results that can be obtained by using recently appeared semantic Web search engines such as semaGER (www.semager.com) and cognition (www.cognition.com). These search engines find semantic keywords and web pages suitable on a semantic base to the context of the search terms specified as input by the user. These engine, to the best of our knowledge, are not intended to analyze and make manifest as links (what we do with our approach) the semantic relations (intended as content similarity and correlation) between the different pages of a given Web site.

In the future, we plan to apply the approach on other Web sites, different in size and application domain, and to conduct controlled experiments to investigate the effectiveness of the recovered navigation maps. These experiments will aim at assessing whether the enhanced version of the sites better satisfies the users’ expectations in terms of contents navigation and information access.

Future work will also be devoted to investigate the effect of adopting and combining different page similarity measures to group pages with similar or related content. The possibility of using different pruning thresholds of the edges of the adopted clustering algorithm will be also considered. The effect of using different singular values of the dimensionality reduction of the latent structure of a Web site will be investigated as well. Moreover, since we used the same stop word list for all the considered Web sites it would be useful to investigate how different stop word lists affect the overall quality on the automatically identified semantic navigation maps.

It will be also worth extending both the approach and the tool prototype to make them suitable for dynamic Web sites, i.e., Web sites for which the content showed into pages and the accessible pages varies depending on some context variable (the user profile, location, etc.) [8][10]. The approach and the tool will be also extended to analyze PDF and Word files. Feature to completely customize the approach (e.g., using personalized stop word list, specifying different criterions for the pruning threshold of the used clustering algorithm) will be also implemented in the new version of the plug-in.

We are currently working on developing software components to be integrated in different and widely employed Web applications, e.g., CMSs, e-learning platforms, and e-commerce application frameworks. With the intent of satisfying the user’s expectations and requirements, we will also investigate the possibility of adapting the navigation maps according to the user’ profile and/or preferences. For example, pages could gain positions (a better score) in the semantic navigation maps in case the user has previously navigated them.

Finally, we are considering the feasibility of extending our approach by implementing a browser plug-in that gives the user the possibility to analyze and produce semantic navigation maps for a Web site of his/her choice.

References

- [1] G. Antoniol, G. Canfora, G. Casazza, and A. De Lucia. "Web Site Reengineering using RMM". *Proc. of the 2nd International Workshop on Web Site Evolution*, Zurich, Switzerland, 2000, pp. 9-16.
- [2] M. Bernardi, G. A. Di Lucca, and D. Distanto, "Reverse Engineering of Web Applications to Abstract User-Centered Conceptual Models". *Proc. of the 10th International Symposium on Web Site Evolution*, IEEE Press, 2008, pp. 55-64.
- [3] C. Boldyreff and P. Tonella. "Web Site Evolution". Special Issue of the *Journal of Software Maintenance*, Vol. 16, No. 1-2, 2004, pp. 1-4.
- [4] C. Boldyreff and R. Kewish. "Reverse Engineering to Achieve Maintainable WWW Sites". *Proc. of the 8th IEEE Working Conference on Reverse Engineering*, Stuttgart, Germany, IEEE CS Press, 2001, pp. 249-257.
- [5] J. Cabot and C. Gómez, "A Catalogue of Refactorings for Navigation Models", *Proc. of the 8th International Conference on Web Engineering*, Yorktown Heights, New York, IEEE CS Press, 2008, pp. 75-85.
- [6] S. Ceri, P. Fraternali, and A. Bongio. "Web Modeling Language (WebML): a Modeling Language for Designing Web Sites". *Computer Networks* 33(1-6): 137-157, 2000.
- [7] D. Cran, E. Pascarello and J. Darren. "Ajax in Action" Manning Publications Co. October, 2005. ISBN: 1932394613
- [8] A. De Lucia, G. Scanniello, and G. Tortora "Identifying Similar Pages in Web Applications using a Competitive Clustering Algorithm". In *Journal on Software Maintenance and Evolution*, Vol. 19, No. 5, September-October 2007, Wiley, pp.: 281-296.
- [9] A. De Lucia, M. Risi, G. Scanniello, and G. Tortora "Clustering Algorithms and Latent Semantic Indexing to Identify Similar Pages in Web Applications", *Proc. of the 9th IEEE International Symposium on Web Site Evolution*, Paris, France, October 5-6, 2007, IEEE CS Press, pp. 65-72.
- [10] A. De Lucia, R. Francese, G. Scanniello, and G. Tortora. "Identifying Cloned Navigational Patterns in Web Applications". In *Journal of Web Engineering* Vol. 5, No. 2, Rinton Press, 2006, pp. 150-174.
- [11] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. "Indexing by Latent Semantic Analysis", *Journal of the American Society for Information Science*, No. 41, 1990, pp. 391-407.
- [12] G. A. Di Lucca, M. Di Penta, and A. R. Fasolino. "An Approach to Identify Duplicated Web Pages". *Proc. of the 26th Annual International Computer Software and Application Conference*, Oxford, UK, IEEE CS Press, 2002, pp. 481-486.

- [13] G. A. Di Lucca, M. Di Penta, G. Antoniol, and G. Casazza. "An Approach for Reverse Engineering of Web-based applications". *Proc. of the 8th IEEE Working Conference on Reverse Engineering*, Stuttgart, Germany, IEEE CS Press, 2001, pp. 231-240.
- [14] D. Distanto, G. Rossi, G. Canfora, and S. Tilley. "A Comprehensive Design Model for Integrating Business Processes in Web Applications". *International Journal of Web Engineering and Technology*, Vol. 2, No. 1, 2007, pp 43-72. Inderscience Publishers, 2007.
- [15] D. Eichmann "Evolving an Engineered Web". *Proc. International Workshop Web Site Evolution*, Atlanta, GA, 1999, pp. 12-16.
- [16] P.J. Flynn, A. K. Jain, and M. N. Murty "Data Clustering: A Review". In *ACM Computing Surveys*, Vol. 31, No. 3, 1999, pp. 264-323.
- [17] A. Garrido, G. Rossi, and D. Distanto, "Model Refactoring in Web Applications". *Proc. of the 9th International Symposium on Web Site Evolution*, IEEE CS Press, 2007, pp. 89-96.
- [18] F. Garzotto and V. Perrone. "On the Acceptability of Conceptual Design Models for Web Applications". In *Proc. of Conceptual Modeling for Novel Application Domains – ER'03 Workshops*. (Chicago, US, Oct.03), LNCS – 2814/ 2003, p. 92-104.
- [19] L. Guttman. "Some necessary conditions for common factor analysis". *Psychometrika*, Vol. 19, 1954, pp. 149-61.
- [20] D. Harman. "Ranking Algorithms", In *Information Retrieval: Data Structures and Algorithms*, Prentice-Hall, Englewood Cliffs, NJ, 1992, pp. 363–392.
- [21] H. F. Kaiser. "The Application of Electronic Computers to Factor Analysis". *Educational and Psychological Measurement*, Vol. 20, 1960, pp. 141-51.
- [22] G. Kappel, B. Pröll, S. Reich, W. Retschitzegger (Eds.), "Web Engineering: The Discipline of Systematic Development of Web Applications", Wiley, 2006.
- [23] N. Koch, A. Kraus, and R. Hennicker. "The Authoring Process of the UML-based Web Engineering Approach" *Proc. of the 1st International Workshop on Web-Oriented Software Technology*, Valencia, Spain (2001), 2001, pp. 105-119.
- [24] A. Kuhn, S. Ducasse, and T. Girba. "Enriching Reverse Engineering with Semantic Clustering", *Proc. of 12th Working Conference on Reverse Engineering*, IEEE CS Press, 2005, pp. 10-20.
- [25] T. K. Landauer and S. T. Dumais. "Solution to Plato's Problem: The Latent Semantic Analysis Theory of Acquisition, Induction and Representation of Knowledge" *Psychological Review*, 1997, Vol. 104, No. 2, pp. 211-240.
- [26] V. L. Levenshtein, "Binary Codes Capable of Correcting Deletions, Insertions, and Reversals". *Cybernetics and Control Theory*, Vol. 10, 1966, pp. 707-710.
- [27] D. Lowe and X. Kong, "NavOptim Coding: Supporting Website Navigation Optimisation using Effort Minimisation". In *2004 IEEE/WIC/ACM International Conference on Web Intelligence*, Beijing, China, 2004, IEEE CS Press, pp. 91-97.
- [28] J. I. Maletic and A. Marcus, "Supporting Program Comprehension Using Semantic and Structural Information". In *Proceedings of 23rd International Conference on Software Engineering*, Toronto, Ont., Canada, 2001, pp. 103-112.
- [29] P. Nakov. "Latent Semantic Analysis for German Literature Investigation" *Proc. of the International Conference, 7th Fuzzy Days on Computational Intelligence, Theory and Applications*, London, UK, 2001, Springer-Verlag, pp. 834–841.

- [30] A. M. Oudshoff, I.E. Bosloper, T. B. Klos, and L. Spaanenburg, "Knowledge Discovery in Virtual Community Texts: Clustering Virtual Communities". In *Journal of Intelligent and Fuzzy Systems*, Vol. 14, No. 1, 2003, pp. 13-24.
- [31] J.M. Pearson, and A. Pearson, "An Exploratory Study into Determining the Relative Importance of Key Criteria in Web Usability: A Multi-Criteria Approach". In *Journal of Computer Information Systems*, July 2008.
- [32] F. Ricca and P. Tonella, "Understanding and Restructuring Web Sites with ReWeb", *IEEE Multimedia*, Vol. 8, No. 2, 2001, pp. 40-51.
- [33] F. Ricca and P. Tonella, "Using Clustering to Support the Migration from Static to Dynamic Web Pages". *Proc. of International Workshop on Program Comprehension*, Portland, Oregon, USA, 2003, pp. 207-216.
- [34] F. Ricca, P. Tonella, C. Girardi, and E. Pianta, "Improving Web Site Understanding With Keyword-Based Clustering" In *Journal of Software Maintenance and Evolution: Research and Practice*, Vol. 20, No. 1, 2008, pp. 1-29.
- [35] D. Schwabe and G. Rossi, "An Object-Oriented Approach to Web-Based Application Design". *Theory and Practice of Object Systems (TAPOS)*, Special Issue on the Internet, Vol. 4, No. 4, October, 1998, pp. 207-225.
- [36] G. Scanniello, D. Distanto, and M. Risi, "Using Semantic Clustering To Enhance the Navigation Structure of Web Sites". *Proc. of the 10th International Symposium on Web Site Evolution*, IEEE CS Press, 2008, pp. 55-64.
- [37] P. Tonella, F. Ricca, E. Pianta, and C. Girardi, "Restructuring Multilingual Web Sites". *Proc. of the 18th International Conference on Software Maintenance (ICSM 2002)*, Montreal, Canada, IEEE CS Press, 2002, pp. 290-299.
- [38] F. Wild, C. Stahl, G. Stermsek, G. Neumann, and Y. Penya. "Parameters Driving Effectiveness of Automated Essay Scoring with LSA". *Proc. of the 9th Computer Assisted Assessment Conference (CAA 2005)*, Loughborough, UK, pp.485-494.
- [39] C. Wohlin, P. Runeson, M. Host, M. C. Ohlsson, B. Regnell, and A. Wesslen, "Experimentation in Software Engineering - An Introduction", Kluwer Academic Publishers Group, 2000.
- [40] G. Tsakonas, C. Papatheodorou, "Exploring Usefulness and Usability in the Evaluation of Open Access Digital Libraries". In *International Journal of Information Processing and Management*, Vol. 44, No. 3, pp. 1234-1250. Pergamon Press, Inc., 2008.
- [41] S. Tilley, "Ten years of Web Site Evolution". *Proc. of the 10th IEEE International Symposium on Web Site Evolution*. IEEE Press, 2008, pp. 11 - 17.