



HAL
open science

Data Mining Special Issue: Mining relevant information on the web: a clique-based approach

Massimiliano Caramia, Giovanni Felici

► To cite this version:

Massimiliano Caramia, Giovanni Felici. Data Mining Special Issue: Mining relevant information on the web: a clique-based approach. *International Journal of Production Research*, 2006, 44 (14), pp.2771-2787. 10.1080/00207540600693713 . hal-00512906

HAL Id: hal-00512906

<https://hal.science/hal-00512906>

Submitted on 1 Sep 2010

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

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



Data Mining Special Issue: Mining relevant information on the web: a clique-based approach

Journal:	<i>International Journal of Production Research</i>
Manuscript ID:	TPRS-2006-IJPR-0194
Manuscript Type:	Original Manuscript
Date Submitted by the Author:	08-Mar-2006
Complete List of Authors:	Caramia, Massimiliano; CNR, IAC Felici, Giovanni; CNR, IASI
Keywords:	WEB-BASED MANUFACTURE, DATA MINING, COMBINATORIAL OPTIMIZATION, DECISION SUPPORT SYSTEMS, E-MANUFACTURING, E-BUSINESS
Keywords (user):	
<p>Note: The following files were submitted by the author for peer review, but cannot be converted to PDF. You must view these files (e.g. movies) online.</p>	
<p>tprs-2006-ijpr-0194final.tex</p>	



Mining Relevant Information on the Web: A Clique-Based Approach

Massimiliano Caramia * Giovanni Felici †

Abstract

The role of information management and retrieval in production processes has been gaining importance in recent years. In this context, the ability to search and quickly find the little information needed in the huge amount of information available has crucial importance. One category of tools devoted to such a task is represented by search engines. The satisfaction of the basic needs of the web user led to the research of new tools that aim at helping more sophisticated users (communities, companies, interest groups) with more elaborated methods. An example is the use of clustering and classification algorithms or other specific data mining techniques. In such context, a proper use of a thematic search engine is a crucial tool to support and orient many activities. Several practical and theoretical problems arise in developing such tools, and we try to face some of them in this paper, extending some previous work on web mining [6]. Here we consider two related problems: how to select an appropriate set of keywords for a thematic engine taking into account the semantic and linguistic extensions of the search context, and how to select and rank a subset of relevant pages given a set of search keywords. Both problems are solved using the same framework, based on a graph representation of the available information and on the search of particular node subsets of such graph. Such subsets are effectively identified by a maximum-weight clique algorithm customized ad-hoc for the specific problems. The methods have been developed in the framework of a funded research project for the development of new web search tools, they have been tested on real data, and are currently being implemented in a prototypal thematic search engine. The web mining method presented in the paper can be applied to web-based design and manufacturing.

*Institute of Applied Mathematics “M. Picone” IAC-CNR, Viale del Policlinico, 137 - 00161 Rome, Italy. E-mail: m.caramia@iac.cnr.it

†Institute of System Analysis and Computer Science “A. Ruberti” IASI-CNR, Viale Manzoni, 30 - 00185 Rome, Italy. E-mail: felici@iasi.cnr.it

1
2
3 **Keywords** Web search, page vectorization, thematic search, maximum clique,
4 *k*-densest subgraph.
5
6

7 8 **1 Introduction** 9

10
11 The importance of information management in production processes has been steadily
12 growing with the evolution of information and communication technology and with the
13 recent globalization process. The consolidated presence of information management has
14 been joined with emerging importance by information retrieval, that is, by the technolo-
15 gies that enable their users to gain information about the production process, market
16 competitors, and customers.
17

18
19 Today the web represents a key driving force for a large spectrum of applications in
20 which users interact with or within companies, organizations, governmental agencies, and
21 educational or collaborative environments. In particular, the area of collaborative en-
22 terprise networks, where different actors share through the web - directly or indirectly -
23 information about their production processes, seems particularly interesting for the appli-
24 cation of thematic search and for the solution of the related problems that arise in this
25 context.
26

27
28 User preferences and expectations, together with usage, content, and structural pat-
29 terns obtained from the web, form the basis for intelligent, personalized, and business-
30 optimal services. Key Web business metrics enabled by proper data collection and pro-
31 cessing are essential to run an effective business or service. Enabling technologies include
32 data mining, scalable data warehousing and preprocessing, sequence discovery, real time
33 processing, document classification, user modeling and evaluation models. Recipient tech-
34 nologies that demand user profiling and usage patterns include recommendation systems,
35 Web analytics applications, content management systems, and fraud or intrusion detection
36 systems.
37

38
39 Another application is web mining for e-commerce. Typical concerns in e-commerce
40 include improved cross-sells, up-sells, personalized ads, targeted assortments, improved
41 conversion rates, and measurements of the affectivity of actions. Other applications are
42 concerned with recommendation and personalization systems, intelligent web services, con-
43 textual information access and retrieval, alert and information filtering systems, adaptive
44
45
46
47
48

1
2
3 hypertext systems, web mining applications for business and competitive intelligence, log
4 analysis for security applications.

5
6 For all the tasks described above it is also very important to obtain a general consensus,
7 among the specific group of users, on what makes some piece of information relevant, and
8 on the main definitions and meaning of the words used within that context; a specific way
9 to tackle these problems is found in the recent development of ontologies (methods to
10 build relations among the meanings of context-oriented terms).

11
12 In this framework, several have considered the use of data mining and optimization
13 techniques, often referred to as *web mining* (for a recent bibliography on this topic see,
14 e.g., [14, 23]). Such techniques are needed to support sophisticated applications, as the
15 standard search methods used in classical search engines are designed for a too general
16 category of users.

17
18 In previous work [6] we proposed a method for improving standard search results in a
19 thematic search engine, where the documents and the pages made available are restricted
20 to a finite number of topics, and the users are considered to belong to a finite number of
21 user profiles. The method uses clustering techniques to identify, in the set of pages resulting
22 from a simple query, subsets that are heterogeneous with respect to a vectorization based
23 on context or profile; small and potentially good subsets of pages are constructed extracting
24 from each cluster the pages with higher scores. Operating on these subsets with a genetic
25 algorithm, the subsets with a good overall score and a high internal dissimilarity are
26 identified. This provides the user with a few non-duplicated pages that represent, from a
27 different viewpoint, the structure of the initial set of pages.

28
29 In this paper we consider two data mining problems related to the classification and
30 search procedures that take place inside a thematic search engine when the pages must be
31 stored and retrieved to satisfy a user search request. The solution of these problems would
32 enhance the capability of a search engine, for example to identify structural information
33 in a set of pages, or to provide semantic extension to search keywords.

34
35 Although such problems may be found in many application environments, we believe
36 that they can be applied successfully to the many web mining problems mentioned above
37 that arise in production processes. The quality of the information retrieved on the web
38 is indeed important when collaborative networks among enterprises that operate in the
39 same logistic chain are set up, or when information on a specific market is retrieved from
40

1
2
3 the web.

4
5 The work presented in this paper refers to some experimental work that has been
6 carried out in a research project (see also [6]) with the aim of developing a thematic
7 engine with new search and information retrieval capabilities to provide more quality
8 when focused on specific topic, category of users, or business target.
9

10
11 Web mining methods such as those presented in this paper find a novel and promising
12 application in web-based design and manufacturing.
13

14
15 The paper is organized as follows. In Section 2 we describe the current architecture
16 of the thematic search engine that is being developed and highlight its features, followed
17 by the notation that will be used throughout the paper (Section 3). Such considerations
18 bring to the surface some of the current limits that we try to overcome with the techniques
19 presented. Here we identify two main problems: the identification of a proper *vectorization*
20 strategy that can be used to represent the pages in a vector space where data mining tech-
21 niques may be used for advanced search; and the identification of semantically connected
22 subsets of pages. Solving the former problem is of crucial importance for the effective-
23 ness of the search engine, as it is the first step of information compression and affects the
24 quality of most of the features of the engine. A solution to the latter problem intends to
25 overcome the typical drawbacks that a user encounters when he/she is submitted a set of
26 pages selected according to some score function that is added on the single pages of the
27 set rather than computed on the sets as a whole.
28

29
30 These two problems, and the models proposed to represent them, are the topic of
31 Sections 4 and 5, respectively. In these two sections we show how the main computational
32 step to solve the two problems is the identification of subgraphs of a given graph G with
33 certain properties, namely some variants of two strongly related problems: the maximum-
34 weight clique and the k -densest subgraph. Given the particular nature of these problems
35 in our application we propose an original solution algorithm that is able to guarantee good
36 performances in very little time and can find the optimal solution to the problem when
37 the graph is of reasonable dimension (Section 6). Computational experience on real data
38 derived from the prototype of the search engine is reported in Section 7, where we describe
39 the behavior of the subset selection algorithm and indirectly measure the effects of the
40 proposed method on the pages considered. Finally, some conclusions are drawn in Section
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
8.

2 Architecture of the Thematic Search Engine

The search engine considered in this paper is a prototype currently being developed by an Italian company in collaboration with IASI-CNR. Its scope is to meet the demand of specific communities of private or professional users (associations, interest groups, large companies) providing customized services. Such services are mainly based on data mining techniques.

The web pages are collected from the web using some ad-hoc spiders that start their search from a set of sensible keywords (context keywords) that are provided by the designer or extracted from a corpus of pages and/or documents.

A given customization of the engine will be composed by a set of *contexts*, that identify some particular topic or category of knowledge, and a set of *profiles*, that identify groups of users that deploy the search engine with homogeneous needs. Contexts and profiles are identified by a limited set of relevant keywords (in the order of few hundreds) that are dynamically evolved according to the behavior of the user and to the pages that are collected from the web. Whenever an advanced search function is executed within a context or a profile, the pages are vectorized according to a transformation of the frequencies of the keywords of that context/profile and data mining is applied; in particular,

- page clustering and genetic algorithms are used for the identification of heterogeneous subsets of pages matching a particular search,
- tree structures are used to enable the user to navigate a set of pages with additional perspectives,
- neural networks are trained on line to learn the preferences of the user.

The main aim of the advanced search functions is to identify meaningful subsets of pages that are relevant both for their matching with some user provided search keywords and for their intrinsic structure, i.e., the relations among the selected pages. One of the additional features that are being implemented in the engine is the construction of a similarity measure among words that is used to extend the search in a semantic sense.

As anticipated, we consider here two particular problems:

- how to select a compact subset of keywords from a larger set to perform the context/profile vectorization in an effective way that may exploit the similarities among

1
2
3 words;

- 4
5
6
7
8
9
- how to select a subset of pages that, beside being matched with some search keywords, form a strongly connected cluster with respect to a context or a profile.

10 In the following sections an approach to deal with these problems is described.

11 12 13 3 Notation

14
15
16 In this section we provide some notation. We are considering the data stored in the search
17 engine after the page collection phase (spidering); such data may be synthesized by the
18 pages, the words contained in these pages, and the frequency of each word in each page.

19
20 We will thus have:

- 21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
- P : set of pages, page $p \in P$;
 - W : set of words, word $w \in W$;
 - $M_{W,P}$: occurrence matrix of words in W in pages in P , where the cell M_{w_j,p_i} represents the number of times word $w_j \in W$ appears in page $p_i \in P$;
 - $i_W(p_1, p_2)$: number of words appearing both in pages p_1 and p_2 (*word-intersection* of p_1, p_2).
 - $i_P(w_1, w_2)$: number of pages containing both the words w_1 and w_2 (*page-intersection* of w_1, w_2).

42 A similarity function among words is also introduced. This function is intended to take
43 into account one or more aspects of word similarity from three standpoints: grammatical,
44 semantic, and linguistic. The first (grammatical) is dealt with using standard language-
45 specific stemming algorithms; it measures whether two words derive from the same root.
46 The second (semantic) takes into account similar meaning of words, and is based on the
47 information contained in thesauri and lists of synonyms. Such lists are built ad-hoc for each
48 specific context/profile of the search engine. The third similarity (linguistic) is verified
49 when the meaning of the two words match in one of the supported languages. In this
50 context we do not enter the details of how these measures are computed, but we assume
51 that some overall similarity measure among words is defined according to the scope of the

specific search, where a user may request to measure similarity by only one of the three measures above or by all of them, disjunctively. Thus, we define:

- $s_G(w_1, w_2)$ grammatical similarity (stemming);
- $s_S(w_1, w_2)$ semantic similarity (synonym);
- $s_L(w_1, w_2)$ language similarity (translation);
- $s(w_1, w_2)$ overall similarity (a combination of the selected above).

According to the specific application, different ways of measuring the similarity among words may be adopted. In the simplest setting, we assume similarity to vary between 1 (maximum similarity) and 0 (no similarity), and state that we use $s_G(w_1, w_2) = 1$ only when w_1 and w_2 come from the same root, $s_S(w_1, w_2) = 1$ only when w_1 and w_2 are synonyms, and $s_L(w_1, w_2) = 1$ only when w_1 and w_2 have the same meaning in different languages.

We can now introduce the models for the two problems described above.

4 The Identification of the Vectorization Set

Data mining algorithms are generally applied on a matrix representation of the data, where the rows are associated with the observed *records* and the columns are associated with the *variables* measured on the records. In the case of web search, the records are typically the pages contained in the engine data base, while the variables that describe them are some measure of their characteristics, derived from word frequencies and page structure (e.g., html tags, meta tags, and other structural information). It is of crucial importance for the effectiveness of the algorithms and for the compactness of the stored data that the *dimension* of the space where the pages are represented (that is, the number of columns of the data matrix) is limited. In previous work [6] we have shown a good experimental performance of clustering and genetic algorithms on a limited (50 – 100) number of “special” words that are used to represent the pages within a context/profile.

The problem that we formulate here is how to choose this limited subset to optimize the retained information. Assume thus that a given set of pages P can be associated with a subset of words W (e.g., all significant words that appear in pages of P at least

1
2
3 k times, or the words that are considered relevant within a context by experts). Our
4 intention is to select a subset W^* of given size with the objective of maximizing the
5 information maintained by W^* , where $|W^*| \ll |W|$. Similar problems are faced by most
6 *feature selection* techniques, where greedy selection techniques are typically used when
7 many variables are present.
8
9

10
11 Consider a subset of words \bar{W} and a set of pages P , and the related occurrence matrix
12 $M_{\bar{W},P}$, whose cell $M_{i,j}$ represents the number of times word w_j appears in the page p_i .
13 A direct measure of the information contained in this matrix is the density of $M_{\bar{W},P}$,
14 $d(M_{\bar{W},P})$, that is equal to the ratio between the number of non-zero cells of the matrix
15 and the total number of cells.
16
17

18
19 If we fix a value h for the dimension of set W^* , the problem is to determine a subset
20 W^* such that $|W^*| = h$ and $d(M_{W^*,P})$ is as large as possible. Moreover, we would also
21 like the words in set W^* to be sufficiently different among each other, in order to enhance
22 their ability to cover different pages. Such objective may be represented as described
23 below. Assume that set W^* is used to generate an *extended occurrence matrix* $M_{W^*,P}^+$
24 with the following procedure. For each pair of words $w_i \in W^*$ and $w_j \in W$, define
25 $\delta_{ij} = 1$ if the similarity between words w_i and w_j is below a given threshold, and $\delta_{ij} = 0$
26 otherwise. Now, define the generic cell associated to keyword w_i and page p_k as $M_{w_i,p_k}^+ =$
27 $\sum_{w_j \in W} \delta_{i,j} \times M_{w_j,p_k}$. Matrix $M_{W^*,P}^+$ is thus determined by an expansion of the occurrence
28 counts taking into account not only the occurrences of the words in the “special” set
29 W^* , but also those of words that are *similar* to the words in W^* . Obviously it is of
30 great interest to obtain an extended occurrence matrix of maximum density, such that the
31 page vectorization, based on similarity rather than on exact matching, results in as much
32 information as possible.
33
34
35
36
37
38
39
40
41
42
43
44

45 The maximization of the retained information in this setting cannot be formulated and
46 solved in a straightforward way, due to the highly combinatorial nature of the problem, to
47 the dimensions of the instances that arise in real applications, and to the particular type
48 of objective function. We thus propose a model based on a graph representation of the
49 problem, where the solution is identified by a subset of the nodes that maximizes a linear
50 combination of the weights on the nodes and on the arcs. The model has two, possibly
51 competing, objectives:
52
53
54
55
56

- 57 1. To maximize the number of pages covered, we want to choose words that appear in

many pages;

2. To avoid the selection of similar words and thus maximize the extended occurrence matrix, we want the sum of the similarity among the words in W^* to be minimal.

Consider now the following complete graph $G = (V_W, E)$, where a node v is present for each word $w \in W$ with weight $a_v = \sum_{p \in P} M_{w,p}$, and an edge is present between each pair of nodes (u, v) with weight a_{uv} proportional to the inverse of the similarity between the words associated with u and v . A good way to represent the two objectives described above would be to identify a subset of h nodes that maximizes the sum of the weights of the selected nodes (page coverage) and the sum of the weights of the edges among the selected nodes (words dissimilarity). More formally, we can state the problem as follows:

$$\begin{aligned} \max \quad & \alpha \cdot \sum_{v \in S} a_v + \beta \cdot \sum_{v \in S} \sum_{w \in S} a_{v,w} \\ & S \subset V_W \\ & |S| = h \end{aligned}$$

where α and β are the overall weights associated with the nodes and edges components of the objective function, respectively. The problem as stated above is a particular type of subgraph selection problem, for which we adopt a solution approach tailored ad-hoc, described in Section 6 and then applied for experimental tests (Section 7.2).

We now turn to the second web mining problem discussed in this paper, the search of subsets of “semantically connected” pages.

5 The Search of Semantically Connected Pages

When a user submits a query to a search engine he/she provides one or more keywords and is returned a set of pages that contain these keywords; the pages are ranked according to some criteria, usually based on the number and type of occurrences of the keywords in the pages. The ranking function differs from search engine to search engine, and strongly characterizes the quality of the engine itself. One problem in ranking is that it is a function that evaluates each page independently, and does not take into account a set of pages as a whole. Such aspect has some relevance in the quality of the search, as the

1
2
3 user normally visualizes only the top pages of the ranking assuming that they are the
4 most interesting from his/her perspective. In [6], this problem has been considered and a
5 method to determine a small subset of pages compliant with the scope of the search, but
6 as different as possible from each other, has been proposed.
7
8
9

10 Here we consider the same problem from a different perspective: we wish to select a
11 subset of pages that are compliant with the search keywords, but that also are characterized
12 by the fact that they share a large subset of words different from the search keywords.
13 We believe this characteristic represents a sort of *semantic connection* of these pages that
14 may be of use to spot some particular aspect of the information present in the pages. In
15 addition, we propose to represent such characteristic of a set of pages by some particular
16 subgraph property that can be sought for by a proper algorithm, similarly to what is done
17 in Section 4.
18
19
20
21
22
23

24 Suppose a search on a set of pages P is required based on the search keywords $SW =$
25 $(sw_1, sw_2, \dots, sw_k)$, and that the pages are vectorized according to a set of relevant words
26 W . Given a set of words $W^o \in W$, let P_{W^o} be the set of pages *induced* by W^o , that is,
27 the pages that contain at least one of the words in W^o ; more formally,
28
29
30
31
32
33

$$P_{W^o} = \{p \in P | \exists w \in W^o : M_{w,p} > 0\}$$

34 Consider now two words w_1 and w_2 in W . We want to express a measure of how these
35 two words are connected with respect to the pages where they appear; naturally, such
36 measure may be given by their *page-intersection* $i_P(w_1, w_2)$; if we restrict the pages used
37 to compute the page-intersection to the pages that contain at least one search keyword,
38 we get a measure of the connection strongly related with our objective: words highly
39 connected are words that co-appear in pages that are relevant for the search.
40
41
42
43
44

45 The idea is thus to select pages that, beside being relevant for the search, are induced
46 by a set of strongly connected words. Moreover, we are inclined to select in this set those
47 words that have a high degree of similarity with the search keywords, to enhance the
48 significance of the pages induced.
49
50
51

52 We can formulate this problem as a *maximum-weight clique* problem on the following
53 graph $G(\lambda) = (V_W, E)$, where:
54
55
56

- 57 • a node v is present for each word $w \in W \setminus SW$, with weight $a_v = \sum_{sw \in SW} s(w, sw)$,

1
2
3 that measures the similarity of the node/word with the search keywords;

- 4
5
6 • an edge is present between nodes u and v representing words w_1 and w_2 if the page-
7 intersection $i_{P_{SW}}(w_1, w_2)$ computed on the set of pages P_{SW} is positive and if the
8 overall page-intersection $i_P(w_1, w_2)$ is greater than λ .
9

10
11 The node weight of a clique of $G(\lambda) = (V_W, E)$ measures how much the words in the
12 clique are similar to the search keywords in SW ; moreover, all pairs of words induce at
13 least λ pages, and at least one page containing the search keywords. Higher values of
14 the parameter λ make the graph more sparse and the clique of maximum weight more
15 connected from the semantic standpoint. We are again faced with a hard optimization
16 problem, as we will discuss in Section 6; moreover, in this case we have stronger time
17 constraints on the computation of the solution as this problem may need to be solved
18 while the user is doing the search.
19

20 We adopt the same algorithmic approach used for the identification of the vectorization
21 set (presented in Section 4) - properly adapted - that provides good heuristic solutions
22 when truncated before termination. Experiments are presented and commented in Section
23 7 for different values of $G(\lambda)$.
24
25

26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

6 Maximum-Weight Clique and Subgraph Selection

Maximum-weight clique and subgraph selection problems have been studied under different
angles in the graph theoretical and mathematical programming community, and have
been used in many contexts to model real life problems. We briefly consider below this
family of combinatorial problems providing the basic definitions and some insights in their
complexity and solution algorithms.

Maximum-weight clique problem. A clique C in a graph G is a node induced sub-
graph where each couple of nodes is connected by an edge. The maximum-weight clique
problem consists in finding a clique of largest weight in a graph. Such a problem has been
shown to be NP-hard [12]. Moreover, for every $\delta > 0$, no $(1 + \delta)$ -approximation to the
maximum clique problem exists (unless $P=NP$) and there exists a $\gamma > 0$ such that no
 n^γ -approximation exists, where n is the number of vertices in the graph (unless $P=NP$).
In particular, the current-best approximation ratio for maximum clique that is achievable

1
2
3 in polynomial time is $O(n/\log^2 n)$ [5]. Recent years have seen much progress in under-
4 standing the inapproximability of the problem, culminating in the result of [19] that it
5 cannot be approximated in ZPP to within $n^{1-\epsilon}$ for any fixed $\epsilon > 0$, unless $\text{NP} \subseteq \text{ZPP}$. The
6 reader is referred to [18] for a survey of the (in)approximability of maximum clique.
7
8
9

10 In the literature, many heuristic approaches for the maximum clique problem are based
11 on sequential greedy heuristics [2, 4, 10, 17, 21]. The idea is to build maximal cliques,
12 starting from an empty clique, and iterating through the repeated addition of vertices.
13 To decide which vertex is added, one uses a greedy heuristic such as choosing the vertex
14 that has the highest degree among candidate vertices. To avoid usual greedy traps, greedy
15 heuristics can be improved by injecting a mild amount of randomization combined with
16 multiple restarts. Also, weights used by the greedy heuristic may be adapted from restart
17 to restart as proposed, e.g., in [17, 21]. To improve the quality of a constructed clique, local
18 search can be used to explore its neighborhood, i.e., the set of cliques that can be obtained
19 by removing and/or adding a given number of vertices: local search iteratively moves in
20 the search space composed of all cliques, from a clique to one of its (best) neighbors. To
21 avoid being trapped in local optima, where all neighbors are cliques of smaller sizes, local
22 search may be combined with some advanced meta-heuristics. For example in [1, 20],
23 Simulated Annealing is used to jump out of local optima by allowing moves towards
24 smaller cliques with a probability proportional to a decreasing temperature parameter. In
25 [11, 13, 15], Tabu Search is used to prevent local search from cycling through a small set
26 of good but suboptimal cliques by keeping track, in a tabu list, of forbidden moves among
27 cliques. In [4], Reactive Search enhances Tabu Search by reactively adapting the size of
28 the tabu list with respect to the need of diversification. In [22], local search is combined
29 with a Genetic Algorithm that allows escape from local maxima by applying crossover and
30 mutation operators to a population of maximal cliques.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 **Subgraph selection problem.** The subgraph selection problem is as follows: given
49 a graph $G = (V, E)$ and weights $w_{ij} = w_{ji}$ on the edges $(i, j) \in E$, the lightest (resp.
50 heaviest) k -subgraph problem is solved by determining a subset $S \subseteq V$ with dimension k
51 such that the sum of the weights associated with the edges of the subgraph induced by S
52 on G is minimum (resp. maximum).
53
54
55

56 When unitary weights are associated with the edges, the lightest k -subgraph problem
57
58
59
60

amounts to the selection of the subgraph of dimension k with the minimum number of edges, namely, the sparsest k -subgraph problem. By switching the sense of the objective functions of the two problems above, it is then easy to obtain the heaviest k -subgraph and densest k -subgraph problems. Equivalently, the heaviest k -subgraph and the densest k -subgraph can be obtained from the lightest k -subgraph and the sparsest k -subgraph, respectively, with a trivial transformation of the objective function coefficients.

Feige and Seltser [8, 9] propose a reduction of the densest k -subgraph to the searching of a clique of a certain size in a graph. Moreover, the same authors show with additional arguments that the densest k -subgraph problem is in the NP-complete class even for bipartite graphs with the node degree bounded by 3. Efficient approximation schemes for the densest k -subgraph problem are analyzed; the authors show that, under the condition that G contains a clique on k nodes, an approximation algorithm exists that, for each $0 < \varepsilon < 1$, determines a subgraph on k vertices with at least $(1 - \varepsilon) \binom{k}{2}$ edges with computational complexity equal to $n^{O((1 + \log \frac{n}{k})/\varepsilon)}$.

In [7], a combinatorial feature selection problem is formulated as a *max distinct points* problem that reduces to the problem of the densest k -subgraph. Evidence is provided that an α -approximated algorithm for the max distinct points problem results in a 2α -approximated algorithm for the densest k -subgraph problem, and that an α -approximated algorithm for the *min- l distinct dimension* problem gives an $\alpha(\alpha + 1)$ -approximated algorithm for the densest subgraph problem. On the other hand, the greedy algorithm with backward selection has been shown to provide an approximation ratio of $O(k/n)$ for any k in [3].

The same case has been considered in [16] with a randomized technique that produces an approximation ratio of 0.25. In general, all approximation ratios obtained for this class of problems are constant in k and linear in n .

Analogies between the maximum-weight clique and the subgraph selection problems. When the lightest (heaviest) k -subgraph problem is run on a complete graph, trivially the output subgraph is a clique itself, and thus the problem reduces to finding a clique of given cardinality k with the minimum (maximum) weight of its edges (nodes). Moreover, also if the graph is not a complete one, we can search for a complete subgraph with k nodes with the lightest or heaviest weight. This strong condition on the structure

1
2
3 of the searched subgraph can be useful to enhance the connectivity characteristics of the
4 subset found.
5
6

7 8 **6.1 The Solution Algorithm** 9

10 What we propose in the following is an algorithm to find a k -densest subgraph or a clique
11 in a graph. This algorithm is fast and concentrates on the weight of the cliques of given
12 input size. Thus, differently from the above mentioned algorithms that are concentrated
13 on finding the maximum clique independently of its size, our algorithm focuses on finding
14 cliques of the same requested size selecting the one with maximum weight, an approach to
15 the problem that is particularly suitable for the two subgraph selection problems described
16 in Sections 4 and 5. This characteristic is, in general, not possessed by the maximum clique
17 algorithms mentioned above. Moreover, even if the problem of finding a clique of size k
18 is polynomial with complexity $O(n^k)$, we note that here we do not want a clique of given
19 size but rather a clique of maximum weight among those of a given size.
20
21

22 The algorithm is constructive and works by choosing a node with highest degree (which
23 initializes the set) from which the constructive process starts. The nodes in the neighbor-
24 hood of the starting node are ordered according to non increasing degree, and are stored
25 in a list L . Nodes in L are visited one by one checking if the current clique can be
26 expanded. In case of a positive answer we expand the current clique.
27

28 Once a node is stored in the list, we associate with the node a parameter *distance* that
29 corresponds to the distance of such a node from the starting node; i.e., the distance is 1 if
30 it belongs to the neighborhood of the starting node, the distance is 2 if it belongs to the
31 neighborhood of a node in the neighborhood of the starting nodes, and so on. Moreover,
32 we associate with each node another parameter called *position* that takes into account the
33 position of the node in the ordering of the neighborhood to which it belongs.
34
35

36 Once we run out of nodes in the neighborhood, we consider the neighborhoods of these
37 nodes just visited and these neighborhoods are stored at the end of list L ordered according
38 to their degree as well, each one with its *distance* and *position* values.
39
40

41 Once the storage is ended, the process iterates, i.e., the nodes inserted in the list and
42 already not visited are scanned starting from the head of L , checking whether the clique
43 can be extended, and keeping on examining neighborhoods till all the graph is visited.
44
45

46 Parameters *distance* and *position* help in applying a sort of backtracking procedure
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 when we run out of nodes during the exploration of the node vector, or when the target
4 cardinality clique is reached in the case of an instance where k is assigned (as it is for the
5 problem described in Section 4). In fact, when this happens, we progressively delete all
6 the nodes in the vector associated with the same distance, i.e., having the same parameter
7 *distance*, as the last nodes put in the clique, and restart the search process by the node
8 with the successive *position* within this distance. If there is no node left within that
9 distance, then we delete all the nodes at $distance - 1$, i.e., with the immediate shorter
10 distance, and start the process with a node with the successive value of *position*.

11
12 The described process is repeated for a certain number of iterations; when this limit
13 is reached, the process restarts from scratch with a starting node of the successive highest
14 degree. In the following we give a compact description of the algorithm for the run with
15 the first node with the highest degree.

16 **Scheme of the Algorithm**

- 17 1. Choose a node with highest degree, and initialize the current clique with this node;
18 initialize a list $L = \emptyset$;
- 19 2. Put neighborhood nodes in L according to non increasing order of degree; set
20 $position = 1$ and $distance = 1$;
- 21 3. Select the first node in L already not visited and check if the current clique can be
22 extended with this node; assign to this node parameters *position* and *distance*. In
23 case of a positive answer, extend the current clique. If a target clique size is reached
24 then goto 6;
- 25 4. If L is empty then goto 5, otherwise $position = position + 1$ and goto 3;
- 26 5. If all the nodes in the graph have been visited, goto 6; otherwise, select nodes with
27 $distance = distance + 1$, and put them in L according to non increasing order of
28 degree; set $position = 1$; goto 3;
- 29 6. Update the best clique if the one found is of better weight and apply backtracking:
30 delete from L all the nodes with *distance* equal to that of the last node inserted in
31 the current clique; delete all the nodes in L with such *distance*, and *position* less
32 than or equal to *position* of the last node inserted in the current clique; goto 4.

Dataset Name	Number of Pages	Number of keywords	Number of nonzero elements in occurrence matrix
<i>A</i>	2,544	3,012	4,399
<i>B</i>	3,405	6,443	177,797
<i>C</i>	8,086	10,453	44,619

Table 1: Dimensions of the dataset used for the experiments

7 Experimental Results

The tests of the proposed methods should aim at two objectives. The first is to verify the effectiveness of our algorithms on the formulated problems, with particular attention to the fact that the dimension of the problems involved may be large and we are bound to use heuristics or truncated search; the second is to check the relevance of the proposed methods in improving the quality of a search engine. Unfortunately, the second objective is more difficult to achieve, as it involves complex user evaluations that must follow the development of a proper prototype. Here we consider mainly the first objective, and only partially the second one, that is presently still under test within the project that sponsored this research activity.

The data used for the experiments were extracted from a database of approximately 100,000 pages associated with different contexts collected and vectorized by the search engine since the beginning of 2004. We considered three sub-contexts of different sizes, from which we derive a set of words and the related occurrence matrices. Table 1 reports the sizes of the three datasets, referred to as *A*, *B*, and *C*.

We have tested both methods on these datasets, analyzing the behavior of the algorithms according to their main parameters (Sections 7.1 and 7.2). For Vectorization Set, we considered a dimension ranging from 20 to 120 for datasets *A*, *B*, and *C*. For Page Selection, we also need to consider the set of search keywords that are used to determine the weight of the nodes; here we focus on the sets *A* and *B* and analyze the behavior of the method with 3 and 5 search keywords. In these experiments the graph connectivity was controlled by the threshold on the minimum weights of the arcs to be included in the graph, that ranged from 1 to 7.

7.1 The Vectorization Set

Some of the experiments conducted on the selection of an optimal vectorization set are synthesized in Table 2. As mentioned above, we have considered different dimensions for the target set h (column 1 of Table 2, starting from 20 up to 120 with step 10). Such value is thus the dimension of the clique identified by the algorithm. The objective function to be maximized adopts weights α for its node component and β for its edge component. Obviously, the tuning of these parameters has an impact on the value of the solution, and we have adopted this reasonable strategy that appears to provide sufficiently stable results:

$$- \alpha = \frac{\text{average edge weight}}{\text{average node weight}}$$

$$- \beta = 2 * \frac{\alpha}{h^2}$$

The software required for the experiment was coded in standard C language and run on an Intel 4 2Ghz processor with 1GB RAM under Linux Debian 2.0, compiled with the standard *gcc* options. The maximum running time for each problem was set to 600 seconds, as in this case there is no need for on-line computation. The best subset determined within the time bound was considered as a satisfactory solution.

As anticipated in Section 4, we measure the quality of the target set with the increase in the density of the projected standard and extended occurrence matrices ($d(M_{W,P}), d(M_{W,P}^+)$). Table 2 reports the variation in the density value of the matrices for the target set and the increase with respect to the whole set for the normal (columns 2 and 3 of Table 2) and the extended case (columns 4 and 5 of Table 2).

Despite of the limitation on the running time, the heuristic solutions provided seem to convey very interesting results. The density of the regular matrix is increased by a significant factor for the three datasets when projecting on a set of 50 words; high increments are present also for larger subsets, where the proportion of “good” words may be reduced, and also the quality of the heuristic solution may decrease due to the time bound. Similar behavior is registered when the increase with respect to the extended matrix is considered. We registered also some variability in the results of the subsets, that is likely to be attributed to the different sparsity degree of the occurrence matrices in the three datasets. It is important to note that most data analysis techniques would have

$ W^* $	$d(M_{W^*,P})$	<i>Increase</i>	$d(M_{W^*,P}^+)$	<i>Increase⁺</i>
Dataset A				
20	0.01677	29.42	41.28	13,578.73
40	0.00876	15.37	21.90	7,204.54
60	0.00594	10.42	15.92	5,237.47
80	0.00479	8.40	12.94	4,255.54
100	0.00395	6.93	10.75	3,536.99
120	0.00340	5.96	9.41	3,094.03
Dataset B				
20	0.01475	18.21	13.30	2890.71
40	0.00768	9.48	26.73	5811.24
60	0.00781	9.64	20.68	4496.37
80	0.00526	6.49	17.72	3852.17
100	0.00464	5.73	16.33	3550.12
120	0.00517	6.38	14.55	3162.51
Dataset C				
20	0.0031	4.44	9.23	3036.56
40	0.0035	5.10	11.57	3806.04
60	0.00224	2.93	9.70	3192.90
80	0.00303	4.27	10.66	3499.73
100	0.00249	3.33	8.7	2878.74
120	0.00234	3.08	8.05	2647.52

Table 2: Increase in density for normal and extended occurrence matrices for datasets A, B, and C

1
2
3 problems in managing such large data matrices using some sort of optimality criterion.
4 Moreover, the selected columns also have the feature of being somehow “different” from
5 each other and thus provide a better covering of the pages in the set. It is also interesting
6 to note how the extended density $d(M_{W^*,P}^+)$ (4th column of Table 2) are all significantly
7 larger than one, meaning that the information maintained by the selected columns is very
8 high when we consider also the frequency of the words that are similar to the selected
9 ones.
10
11
12
13
14

15 16 **7.2 Page Selection**

17
18 To investigate the appropriateness of the solution method for the second problem con-
19 sidered (see Section 5) we have run some experiments, described in Tables 3 and 4. The
20 experiments were run on the same platform of the previous case and a maximum running
21 time of 10 seconds was given. In the majority of the cases the best solution was obtained
22 by the proposed heuristic within 3 seconds. Such running times appear reasonable for
23 real time application once the code has been properly optimized and a more up-to-date
24 hardware is used - not to mention the strong parallelization that may be adopted for the
25 clique and subset selection algorithms. The tables give the results of the method when
26 applied to datasets *A* and *B* described above, using two different sets of search keywords
27 with dimension 3 and 5, respectively. For each such set the tables report the value of λ
28 (column 1), that is the minimum value of required page-connectivity of an edge in the
29 graph (see Section 5); the dimension of the best clique determined by the algorithm (col-
30 umn 2); the number of different pages induced by the nodes/words in the clique (column
31 3); the average number of search keywords in the pages, and the average number of words
32 that are similar to the search keywords (columns 4 and 5, respectively), where we declare
33 a word to be similar to a keyword if the similarity measure between the two strings is
34 above a given threshold.
35
36
37
38
39
40
41
42
43
44
45
46
47

48 From the results reported in Tables 3 and 4 we see how the increase in the connectivity
49 threshold λ has in general the effect of reducing the number of selected nodes (words) of
50 the clique, but this increase does not always translates to the number of pages that are
51 induced by the selected words; in fact, as λ increases, the graph becomes sparse and the
52 arcs connect nodes that share more pages. A clique will then be formed by nodes that
53 induce a larger set of pages. The non monotonicity of the relations among the columns of
54
55
56
57
58
59
60

Dataset A			3 Keywords		
Value of λ	Number of nodes	Induced pages	Average number of search keywords	Average number of similar words	
1	18	127	0.91	35.96	
2	13	24	1.05	38.90	
3	14	91	1.19	41.85	
4	14	82	1.32	44.90	
5	12	104	1.04	41.57	
6	11	16	1.64	16.56	
7	11	21	2.24	17.17	

Dataset A			5 Keywords		
Value of λ	Number of nodes	Induced pages	Average number of search keywords	Average number of similar words	
1	12	7	1.7143	10.14	
2	6	45	0.6	6.51	
3	5	45	0.6	6.51	
4	5	45	0.6	6.51	
5	4	45	0.6	6.51	
6	4	45	0.6	6.51	
7	4	13	0.9231	2.69	

Table 3: Selection of pages with semantic connection for dataset A

Dataset B			3 Keywords		
Value of λ	Number of nodes	Induced pages	Average number of search keywords	Average number of similar words	
1	18	27	1.75	47.71	
2	13	35	0.24	29.36	
3	13	32	0.63	29.36	
4	7	29	0.07	29.36	
5	4	42	1.35	29.36	
6	4	38	1.35	29.36	
7	4	17	2.68	40.23	

Dataset B			5 Keywords		
Value of λ	Number of nodes	Induced pages	Average number of search keywords	Average number of similar words	
1	31	97	2.3486	14.36	
2	14	54	4.0556	24.85	
3	14	91	2.5495	13.12	
4	14	82	2.8293	14.56	
5	12	58	2.6923	16.12	
6	11	66	2.8125	16.56	
7	11	70	1.8143	14.71	

Table 4: Selection of pages with semantic connection for dataset B

1
2
3 the table may thus be attributed to the particular characteristics of the data considered.
4 We also note that the last two columns of the tables, that measure the interestingness
5 of the pages with respect to the search keywords, report large numbers and that such
6 numbers increase with the dimension of the search keyword set, as is to be expected.
7
8
9

10 11 8 Conclusions

12
13
14 The work reported in this paper deals with two open problems: the first is related with the
15 representation of the web pages in a search engine, and, more specifically, with a technique
16 to determine small subsets of words that are convenient to use to project pages in a vector
17 space and apply data mining techniques; the second deals with the identification of small
18 sets of pages that have a strong semantic structure, that is, pages that contain many
19 similar words and match one or more user requested keywords. Both methods exploit
20 the potential of the use of a similarity function among words that takes into account,
21 according to a very simple yet effective interpretation of semantics, stemming, synonyms,
22 and translation.
23
24
25
26
27
28
29

30 The two problems have been modeled as particular subgraph selection problems, and
31 a fast heuristic algorithm for their solution has been proposed and tested on real data
32 extracted from a prototypal search engine.
33
34

35 The results provide some interesting insight in the potential of the new problems that
36 can be formulated and efficiently solved to give more value to web search, a very common
37 process in many production processes. Web mining is now frequently applied in web-
38 based design and manufacturing and the use of a thematic approach to search appears
39 particularly relevant in this context.
40
41
42
43

44 The methods described are now being engineered and integrated in the new version
45 of the prototypal search engine to undergo real user impact tests. Future research in this
46 direction will include the tuning of the parameters in the algorithms used to identify the
47 subset of words, and a deeper analysis of the interactions between computational time
48 and solution quality to possibly improve the quality of the method with *ad-hoc* heuristic
49 strategies.
50
51
52
53
54
55
56
57
58
59
60

References

- [1] Aarts, E.H.L., and Korst, J.H.M., *Simulated annealing and Boltzmann machines: a stochastic approach to combinatorial optimization and neural computing*, 1989 (John Wiley & Sons, Chichester, U.K.).
- [2] Abello, J., Pardalos, P.M., and Resende, M.G.C., On maximum clique problems in very large graphs. In *DIMACS Series on Discrete Mathematics and Theoretical Computer Science*, volume 50, pp. 119-130, 2000 (American Mathematical Society).
- [3] Asahiro, Y., Iwama, K., Tamaki, H., and Tokuyama, T., Greedily finding a dense subgraph, in *Proceedings of the 5th Scandinavian Workshop on Algorithm Theory (SWAT), Lecture Notes in Computer Science 1097*, Springer-Verlag, 1996, pp. 136-148.
- [4] Battiti, R., and Protasi, M., Reactive local search for the maximum clique problem. *Algorithmica*, 2001, **29**(4), 610-637.
- [5] Boppana, R.B., and Halld'orsson, M.M., Approximating maximum independent sets by excluding subgraphs. *BIT*, 1992, **32**, 180-196.
- [6] Caramia, M., Felici, G., and Pezzoli, A., Improving Search Results with Data Mining in a Thematic Search Engine. *Computers and Operations Research*, 2004, **31**, 2387-2404.
- [7] Charikar, M., Guruswami, V., Kumar, R., Rajagopalan, S., and Sahai, A., Combinatorial Feature Selection Problems, in *Proceedings of FOCS*, 2000, pp. 631-640.
- [8] Feige, U., and Seltser, M., On the densest k -subgraph problem, Report No. CS97-16, Weizmann Institute of Science, Rehovot, Israel, 1997.
- [9] Feige, U., Kortsarz, G., and Peleg, D., The Dense k -Subgraph Problem. *Algorithmica*, 2001, **29**, 410-421.
- [10] Feo, T.A., Resende, M.G.C., and Smith, S.H., A greedy randomized adaptive search procedure for maximum independent set. *Operations Research*, 1994, **42**, 860-878.
- [11] Friden, C., Hertz, A., and de Werra, D., Stabulus: a technique for finding stable sets in large graphs with tabu search. *Computing*, 1989, **42**, 35-44.

- 1
2
3 [12] Garey, M.R., and Johnson, D.S., *Computers and Intractability*, 1979 (Freeman).
4
5
6 [13] Gendreau, M., Soriano, P., and Salvail, L., Solving the maximum clique problem
7 using a tabu search approach. *Annals of Operations Research*, 1993, **41(4)**, 385-403.
8
9
10 [14] Getoor, L., Senator, T.E., Domingos, P., and Faloutsos, C., *Proceedings of the 9th*
11 *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
12 2003 (Washington DC, USA).
13
14
15 [15] Glover, F., and Laguna, M., Tabu search. In *Modern Heuristics Technics for Combi-*
16 *natorial Problems*, pp. 70-141, 1993 (Blackwell Scientific Publishing, Oxford, UK).
17
18
19 [16] Goemans, M.X., *Mathematical programming and approximation algorithms*, Lecture
20 on Approximate Solution of Hard Combinatorial Problems, 1996 (Summer School,
21 Udine).
22
23
24 [17] Grosso, A., Locatelli, M., and Della Croce, F., Combining swaps and node weights in
25 an adaptive greedy approach for the maximum clique problem. *Journal of Heuristics*,
26 2004, **10(2)**, 135-152.
27
28
29 [18] Halld'orsson, M.M., Approximations of independent sets in graphs, in *Proceedings*
30 *APPROX '98 Conference, Lecture Notes in Computer Science 1444*, Springer-Verlag,
31 1998, pp. 1-13.
32
33
34 [19] Hastad, J., Clique is hard to approximate within $n^{1-\epsilon}$. *Acta Mathematica*, 1999, **182**,
35 105-142.
36
37
38 [20] Homer, S. and Peinado, M., On the performance of polynomial-time clique approx-
39 imation algorithms on very large graphs. In *Cliques, Coloring, 27 and Satisfiability:*
40 *second DIMACS Implementation Challenge*, volume 26, pp. 103-124, 1996 (American
41 Mathematical Society).
42
43
44 [21] Jagota, A., and Sanchis, L.A., Adaptive, restart, randomized greedy heuristics for
45 maximum clique. *Journal of Heuristics*, 2001, **7(6)** pp. 565-585.
46
47
48 [22] Marchiori, E., Genetic, iterated and multistart local search for the maximum clique
49 problem, in *Applications of Evolutionary Computing, Proceedings of EvoWorkshops:*
50
51
52
53
54
55
56
57
58
59
60

1
2
3 *EvoCOP, EvoIASP, EvoSTim, Lecture Notes in Computer Science 2279, Springer-*
4 *Verlag, 2002, pp. 112-121.*
5
6

- 7 [23] Zaiane, O.R., Srivastava, J., Spiliopoulou, M., and Masand, B.M., *Proceedings of the*
8 *International Workshop in Mining Web Data for Discovering Usage Patterns and*
9 *Profiles, 2002 (Edmonton, Canada).*
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only