

Cluster Generation and Cluster Labelling for Web Snippets: A Fast and Accurate Hierarchical Solution

Filippo Geraci

Istituto di Informatica e Telematica, Consiglio Nazionale delle Ricerche
Via G Moruzzi, 1 – 56124 Pisa, Italy and
Dipartimento di Ingegneria dell’Informazione, Università di Siena,
Via Roma, 56 – 53100 Siena, Italy, f.geraci@iit.cnr.it

Marco Pellegrini

Istituto di Informatica e Telematica, Consiglio Nazionale delle Ricerche
Via G Moruzzi, 1 – 56124 Pisa, Italy, m.pellegrini@iit.cnr.it

Marco Maggini

Dipartimento di Ingegneria dell’Informazione, Università di Siena,
Via Roma, 56 – 53100 Siena, Italy, maggini@dii.unisi.it

Fabrizio Sebastiani

Istituto di Scienza e Tecnologia dell’Informazione,
Consiglio Nazionale delle Ricerche
Via G Moruzzi, 1 – 56124 Pisa, Italy, fabrizio.sebastiani@isti.cnr.it

July 13, 2007

Abstract

This paper describes *Armil*¹, a meta-search engine that groups the Web snippets returned by auxiliary search engines into disjoint labelled clusters. The cluster labels generated by *Armil* provide the user with a compact guide to assessing the relevance of each cluster to his/her information need. Striking the right balance between running time and cluster well-formedness was a key point in the design of our system. Both the clustering and the labelling tasks are performed on the fly by processing only the snippets provided by the auxiliary search engines, and use no external sources of knowledge. Clustering is performed by means of a fast version of the furthest-point-first algorithm for metric k -center clustering. Cluster labelling is achieved by combining intra-cluster and inter-cluster term extraction based on a variant of the information gain measure. We have tested the clustering effectiveness of *Armil* against *Vivisimo*, the *de facto* industrial standard in Web snippet clustering, using as benchmark a comprehensive set of snippets obtained from the *Open Directory Project* hierarchy. According to two widely accepted “external” metrics of clustering quality, *Armil* achieves better performance levels by 10%. We also report the results of a thorough user evaluation of both the clustering and the cluster labelling algorithms. On a standard desktop PC (AMD Athlon 1Ghz Clock with 750 Mbytes RAM), *Armil* performs clustering and labelling altogether of up to 200 snippets in less than one second.

¹A Preliminary version of this work has appeared in the proceedings of ACM SAC 2006 [10] and SPIRE 2006 [11]. Work partially supported by the Italian Registry for the ccTLD “.it”.

1 Introduction

An effective search interface is a fundamental component in a Web search engine. In particular, the quality of presentation of the search results often represents one of the main keys to the success of such systems. Most search engines present the results of a user query as a ranked list of Web snippets. Ranking algorithms play a crucial role in this approach, since users usually browse at most the 10 top-ranked items. Snippet quality is also an important issue, since good-quality snippets allow the user to determine whether the referred pages match or not his/her information need. In order to provide a useful hint about the real content of the page, a Web snippet includes both the page title and a short text fragment, that often displays the query terms in context. *Meta-search engines* (MSEs) integrate the items obtained from multiple “auxiliary” search engines, with the purpose of increasing the coverage of the results. However, without an accurate design, MSEs might in principle even worsen the quality of the information access experience, since the user is typically confronted with an even larger set of results (see [27] for a recent survey of challenges and techniques related to building meta-search engines). Thus, key issues to be faced by MSEs concern the exploitation of effective algorithms for merging the ranked lists of results retrieved by the different auxiliary search engines (while at the same time removing the duplicates), and the design of advanced user interfaces based on a structured organization of the results, so as to help the user to focus on the most relevant subset of results. This latter aspect is usually implemented by grouping the results into homogeneous groups by means of clustering or categorization algorithms.

In the context of the the World Wide Web, clustering is not only useful for meta-searching. For example, regular search engine do clustering to some extent to avoid showing the user too many semantically-equivalent documents in the first page of results. This activity is close to the classical *duplicate* or *near-duplicate* detection in information retrieval, and one can take advantage of the fact that duplicates and near duplicate are easy to detect via multiple hashing or shingling techniques, and these tasks can be carried out, at least in part, in an off-line setting (see e.g. [13]). Moreover, in the duplicate detection activity labelling is not an issue.

This paper describes the *Armil* system², a meta-search engine that organizes the Web snippets retrieved from auxiliary search engines into disjoint clusters and automatically constructs a title label for each cluster by using only the text excerpts available in the snippets. Our design efforts were directed towards devising a fast clustering algorithm able to yield good-quality homogeneous groups, and a distillation technique for selecting appropriate and useful labels for the clusters. The speed of the two algorithms was a key issue in our design, since the system must organize the results on the fly, thus minimizing the latency between the issuing of the query and the presentation of the results. Second-level clustering is also performed automatically at query time (i.e. not on demand) to minimize latency. In *Armil*, an equally important role is played by the clustering component and by the labelling component. Clustering is accomplished by means of an improved version of the furthest-point-first (FPF) algorithm for k -center clustering. To the best of our knowledge this algorithm had never been used in the context of Web snippet clustering or text clustering. The generation of the cluster labels is instead accomplished by means of a combination of intra-cluster and inter-cluster term extraction, based on a modified version of the information gain measure. This approach tries to capture the most significant and discriminative words for each cluster.

One key design feature of *Armil* is that it relies almost only on the information returned by the auxiliary search engines, i.e. the snippets; this means that no substantial external source of information, such as ontologies or morphological and syntactic linguistic resources, is used. We use just stop word lists, stemmers and language recognition tools. We thus demonstrate that such a lightweight approach,

²An armillary sphere (also known as a spherical astrolabe, armilla, or *armil*) is a navigation tool in the form of a model of the celestial sphere, invented by Eratosthenes in 255 BC. Renaissance scientists and public figures were often portrayed with one hand on an armil, since it represented the height of wisdom and knowledge (see http://en.wikipedia.org/wiki/Armillary_sphere). The *Armil* query interface can be freely accessed and used at the url <http://armil.iit.cnr.it/>.

together with carefully crafted algorithms, is sufficient to provide a useful and successful clustering-plus-labelling service. Obviously, this assumption relies on the hypothesis that the quality of the snippets returned by the auxiliary search engines is satisfactory.

We have tested the clustering effectiveness of *Armil* against *Vivisimo*, the *de facto* industrial standard in Web snippet clustering, using as benchmark a comprehensive set of snippets obtained from the *Open Directory Project* hierarchy. According to two metrics of clustering quality that are normalized variants of the Entropy and the Mutual Information [4], *Armil* achieves better performance levels by 10%. Note that, since the normalization reduces the ranges of these measures in the interval $[0, 1]$, an increase of 10% is noteworthy.

We have tested the clustering effectiveness of *Armil* against variants of the classical *k*-means clustering algorithm. Our method outperforms the *k*-means clustering algorithm in time-quality trade-offs.

We also report the results of a thorough user evaluation of both the clustering and the cluster labelling algorithms.

1.1 Outline of the clustering algorithm

Clustering and labelling are both essential operations for a Web snippet clustering system. However, each previously proposed such system strikes a different balance between these two aspects. Some systems (e.g. [7, 22]) view label extraction as the primary goal, and clustering as a by-product of the label extraction procedure. Other systems (e.g. [19, 39]) view instead the formation of clusters as the most important step, and the labelling phase is considered as strictly dependent on the clusters found. We have followed this latter approach. In order to cluster the snippets in the returned lists, we map them into a vector space endowed with a distance function, which is a metric; then a modified furthest-point-first algorithm (M-FPF) is applied to generate the clusters. The M-FPF algorithm generates the same clusters of the “standard” FPF algorithm [12, 16], but it uses filters based on the triangular inequality to speed up the computation. As such, M-FPF inherits a very important property of the FPF algorithm, i.e. it is within a factor 2 of the optimal solution for the *k*-center problem [12, 16]. In other words, for any fixed number *k*, M-FPF produces a *k*-clustering (i.e. a partition of the items into *k* non-overlapping clusters) such that the maximum cluster diameter is at most twice as large as that of the “optimal” *k*-clustering (i.e. the one that minimizes such maximum diameter). The competitive property of M-FPF is even stronger: the approximation factor of 2 cannot be improved with any polynomial approximation algorithm, unless $P = NP$ [36]. The strength of this formal property has been our main motivation for selecting M-FPF as the algorithmic backbone for *Armil*. The second interesting property of M-FPF is that it does not compute centroids of clusters. Centroids tend to be dense vectors and, as such, their computation and/or update in high-dimensional space is a computational burden³. M-FPF relies instead only on pairwise distance calculations between snippets, and as such better exploits the sparsity of the snippet vector representations.

1.2 Outline of the cluster labelling algorithm

The cluster labelling phase aims at extracting from the set of snippets assigned to each cluster a sequence of words highly descriptive of the corresponding group of items. The quality of the label depends on its well-formedness (i.e. whether the text is syntactically and semantically plausible), on its descriptive power (i.e. how well it describes what is contained in the cluster), and on its discriminative power (i.e. how well it differentiates what is contained in the cluster with respect to what is

³The clustering literature also discusses the notion of a cluster “medoid”; similarly to a centroid, a cluster medoid plays the role of cluster representative, but typically has a sparse representation. Unfortunately, the methods proposed in the literature for finding high-quality medoids are not compatible with the real-time nature of the envisioned online Web service [42]. In this paper we will give a fast and effective medoid searching technique.

contained in other clusters)⁴. The possibility to extract good labels directly from the available snippets is strongly dependent on their quality and, obviously, on the homogeneity of the produced clusters. In order to pursue a good tradeoff between descriptive and discriminative power, we select candidate words for each cluster by means of IG_m , a modified version of the *Information Gain* measure [4]. For each cluster, IG_m allows the selection of those words that are the most representative of its contents and are the least representative of the contents of the other clusters. Finally, in order to construct plausible labels, rather than simply using the list of the top-scoring words (i.e. the ones that maximize IG_m), the system looks within the titles of the returned Web pages for the substring that best matches the selected top-scoring words.

Once each cluster has been assigned a set of descriptive and discriminative keywords (we call such set of keywords the cluster *signatures*), all the clusters that share the same signature are merged. This reduces the arbitrariness inherent in the choice of their number k , that is fixed *a priori* independently of the query.

1.3 Outline of the paper

The paper is organized as follows. In Section 2 we review related work on techniques for the automatic re-organization of search results. Section 3 introduces the data representation adopted within Armil and sketches the properties of the M-FPF clustering algorithm and of the cluster labelling algorithm. In Section 4 the architecture of Armil is described in detail. The results of the system evaluation are reported in Sections 5, 6,7 and 8. Finally, in Section 9 conclusions and prospective future research are discussed.

2 Previous work

Tools for clustering Web snippets have received attention in the research community. In the past, this approach has had both critics [21, 20] and supporters [35], but the proliferation of commercial Web services such as Copernic, Dogpile, Groxis, iBoogie, Kartoo, Mooter, and Vivisimo seems to confirm the potential validity of the approach. Academic research prototypes are also available, such as Grouper [38, 39], EigenCluster [3], Shoc [41], and SnakeT [7]. Generally, details of the algorithms underlying commercial Web services are not in the public domain.

Maarek et al. [24] give a precise characterization of the challenges inherent in Web snippet clustering, and propose an algorithm based on complete-link hierarchical agglomerative clustering that is quadratic in the number n of snippets. They introduce a technique called “lexical affinity” whereby the co-occurrence of words influences the similarity metric.

Zeng et al. [40] tackle the problem of detecting good cluster names as preliminary to the formation of the clusters, using a supervised learning approach. Note that the methods considered in our paper are instead all unsupervised, thus requiring no labelled data.

The EigenCluster [3], Lingo [28], and Shoc [41] systems all tackle Web snippet clustering by performing a singular value decomposition of the term-document incidence matrix⁵; the problem with this approach is that SVD is extremely time-consuming, hence problematic when applied to a large number of snippets. By testing a number of queries on Eigencluster we have observed that, when operating on many snippets (roughly 400), a reasonable response time (under 1 second) is attained by limiting the number of generated clusters to a number between 5 and 10, and avoiding a clustering decision for over 50% of the data. Zamir and Etzioni [38, 39] propose a Web snippet clustering mechanism (Suffix Tree Clustering – STC) based on suffix arrays, and experimentally compare STC with algorithms

⁴These requirements are akin to some of the desirable properties advocated in [19] for taxonomic trees. In particular property 3 (*Sibling nodes distinctiveness*) and property 4 (*Node label predictiveness*) [19, page 660].

⁵At the time of this writing the Eigencluster system is available on-line at <http://eigencluster.csail.mit.edu/>.

such as k -means, single-pass k -means [25], Backshot and Fractionation [5], and Group Average Hierarchical Agglomerative Clustering (GAHAC). They test the systems on a benchmark obtained by issuing 10 queries to the Metacrawler meta-search engine, retaining the top-ranked 200 snippets for each query, and manually tagging the snippets by relevance to the queries. They then compute the quality of the clustering obtained by the tested systems by ordering the generated clusters according to precision, and by equating the effectiveness of the system with the average precision of the highest-precision clusters that collectively contain 10% of the input documents. This methodology had been advocated in [15], and is based on the assumption that the users will anyway be able to spot the clusters most relevant to their query. Average precision as computed with this method ranges from 0.2 to 0.4 for all the algorithms tested (STC coming out on top in terms of both effectiveness and efficiency). Interestingly, the authors show that very similar results are attained when full documents are used instead of their snippets, thus validating the snippet-based clustering approach.

Strehl et al. [34] experiment with four similarity measures (Cosine similarity, Euclidean distance, Pearson Correlation, extended Jaccard) in combination with five algorithms (k -means, self-organizing maps, random clustering, min-cut hypergraph partitioning, min-cut weighted graph partitioning). Two data sets were used, one formed by 2340 documents pre-classified in 20 categories from the news repository of Yahoo!, the second formed by 966 documents pre-classified in 10 categories from the Business Directory of Yahoo!. Overall quality is measured in terms of normalized mutual information. For a specific clustering also the quality of the single clusters in terms of single cluster purity, single cluster entropy are given. The comparisons are made only in terms of output quality, computing time not being considered relevant in this setting. The highest quality results are obtained via cosine similarity and Jaccard distance combined with min-cut graph partitioning. Our experiments have confirmed that for snippets better results are obtained using variants of Jaccard distance, with respect to standard tf-idf and cosine similarity.

Lawrie and Croft [22] view the clustering/labelling problem as that of generating multilevel summaries of the set of documents (in this case the Web snippets returned by a search engine). The technique is based on first building off-line a statistical model of the background language (e.g. the statistical distribution of words in a large corpus of the English language), and on subsequently extracting “topical terms” from the documents, where “topicality” is measured by the contribution of a term to the Kullback-Leibler divergence score of the document collection relative to the background language. Intuitively, this formula measures how important this term is in measuring the distance of the collection of documents from the distribution of the background language. Additionally, the “predictiveness” of each term is measured. Intuitively, predictiveness measures how close a term appears (within a given window size) to other terms. In the summaries, terms of high topicality and high predictiveness are preferred. The proposed method is shown to be superior (by using the KL-divergence) to a naive summarizer that just selects the terms with highest $tf * idf$ score in the document set.

Kammamuru et al. [19] propose a classification of Web snippet clustering algorithms into *monothetic* (in which the assignment of a snippet to a cluster is based on a single dominant feature) and *polythetic* (in which several features concur in determining the assignment of a snippet to a cluster). The rationale for proposing a monothetic algorithm is that the single discriminating feature is a natural label candidate. The authors propose such an algorithm (a type of greedy cover) in which the snippets are seen as sets of words and the next term is chosen so as to maximize the number of newly covered sets while minimizing the hits with already covered sets. The paper reports empirical evaluations and user studies over two classes of queries, “ambiguous” and “popular”. The users were asked to compare 3 clustering algorithms over the set of queries and, for each query, were asked to answer 6 questions of a rather general nature on the generated hierarchy (not on the single clusters).

Ferragina and Gulli [7] propose a method for hierarchically clustering Web snippets, and produce a hierarchical labelling based on constructing a sequence of la-

belled and weighted bipartite graphs representing the individual snippets on one side and a set of labels (and corresponding clusters) on the other side. Data from the Open Directory Project (ODP)⁶ is used in an off-line and query-independent way to generate predefined weights that are associated on-line to the words of the snippets returned by the queries. Data is collected from 16 search engines as a result of 77 queries chosen for their popularity among Lycos and Google users in 2004. The snippets are then clustered and the labels are manually tagged as relevant or not relevant to the cluster to which they have been associated. The clusters are ordered in terms of their weight, and quality is measured in terms of the number of relevant labels among the first n labels, for $n \in \{3, 5, 7, 10\}$. Note that in this work the emphasis is on the quality of the labels rather than on that of the clusters (although the two concepts are certainly related), and that the ground truth is defined “a posteriori”, after the queries are processed.

3 The clustering algorithm and the labelling algorithm

We now describe in detail the methods used by Arnil for Web snippet clustering and cluster labelling.

3.0.1 The distance function

Each snippet is turned into a “bag of words” after removing stop words and performing stemming. In [10] we report experiments using, as a distance function, (i) the cosine distance measure (i.e. the complement to 1 of the cosine similarity function) applied to vectors of terms weighted by $tf * idf$ ⁷, and (ii) a slight modification of the standard Jaccard Distance, which we call *Weighted Jaccard Distance* (WJD); in those experiments, (ii) has performed at the same level of accuracy as (i), but has proven much faster to compute.

In this journal version we improve on the preliminary results of [10] by using the Generalized Jaccard Distance described in [2]. Given two “bag-of-words” snippet vectors $s_1 = (s_1^1, \dots, s_1^h)$ and $s_2 = (s_2^1, \dots, s_2^h)$, the *Generalized Jaccard Distance* is:

$$D(s_1, s_2) = 1 - \frac{\sum_i \min(s_1^i, s_2^i)}{\sum_i \max(s_1^i, s_2^i)}.$$

The term weights s_a^i consist of “weighted term frequencies”, obtained as weighted sums of the numbers of occurrences of the term in the snippet, where we take advantage of the structure of the snippets by weighting different parts of the snippet (page title, text fragment, URL) differently; more precisely, weight 3 is assigned to a term occurring in the page title, weight 1 to a term occurring in the text fragment, and weight 0 is assigned to a term occurring in the URL (since, in previous experiments we had run, the text of the URL had proven to give no contribution in terms of cluster quality). Note that, when using unit weights only, the Generalized Jaccard Distance coincides with the standard Jaccard Distance. In [2] it is proved that GJD is a metric.

3.1 The clustering algorithm

We approach the problem of clustering Web snippets as that of finding a solution to the classic k -center problem:

Given a set S of points in a metric space M endowed with a metric distance function D , and given a desired number k of resulting clusters, partition S into

⁶<http://www.dmoz.org/>

⁷The acronym $tf * idf$ stands for *term frequency* and *inverse document frequency*. It is a standard technique for assigning real vectors to documents belonging to a set of documents (corpus) based on the number of times a term appears in a document and on the number of documents of the corpus containing that term (see e.g. [1, page 29]).

non-overlapping clusters C_1, \dots, C_k and determine their “centers” $\mu_1, \dots, \mu_k \in M$ so that the radius $\max_j \max_{x \in C_j} D(x, \mu_j)$ of the widest cluster is minimized.

The k -center problem can be solved approximately using the furthest-point-first (FPF) algorithm [12, 16], which we now describe. Given a set S of n points, FPF builds a sequence $T_1 \subset \dots \subset T_k = T$ of k sets of “centers” (with $T_i = \{\mu_1, \dots, \mu_i\} \subset S$) in the following way.

1. At the end of iteration $i - 1$ FPF holds the mapping μ defined for every point $p_j \in S \setminus T_{i-1}$ as:

$$\mu(p_j) = \arg \min_{\mu_s} D(p_j, \mu_s) \quad (1)$$

i.e. the center in T_{i-1} closest to p_j ; $\mu(p_j)$ is called the *leader* of p_j . Note that this mapping is established in the first iteration in time $O(n)$.

2. At iteration i , among all points p_j , FPF picks

$$\mu_i = \arg \max_{p_j} D(p_j, \mu(p_j))$$

i.e. the point for which the distance to its leader is maximum, and makes it a new center, i.e. adds it to T_{i-1} , thus obtaining T_i . This selection costs $O(n)$.

3. Compute the distance of μ_i to any point in $S \setminus T_i$ and update the mapping μ if needed. Thus μ is now correct for the beginning of iteration $i + 1$. This update phase costs $O(n)$.

The final set of centers $T = \{\mu_1, \dots, \mu_k\}$ defines the resulting k -clustering, since each center μ_i implicitly identifies a cluster C_i as the set of data points whose leader is μ_i . Note that T_1 is initialized to contain a single point chosen at random from S ; this random choice is due to the fact that, in practice, both the effectiveness and the efficiency of the algorithm can be seen experimentally to be rather insensitive to this choice.

Most of the computation is actually devoted to computing distances and updating the auxiliary mapping μ : this takes $O(n)$ time per iteration, so the total computational cost of the algorithm is $O(nk)$. Here we define an improved version of this algorithm that exploits the triangular inequality in order to filter out useless distance computations. This modified algorithm (M-FPF), which we now describe, works in any metric space, hence in any vector space⁸.

Consider, in the FPF algorithm, any center $\mu_x \in T_i$ and its associated set of closest points $N(\mu_x) = \{p_j \in S \setminus T_i \mid \mu(p_j) = \mu_x\}$. We store $N(\mu_x)$ as a ranked list, in order of decreasing distance from μ_x . When a new center μ_y is added to T_i , in order to identify its associated set of closest points $N(\mu_y)$ we scan every $N(\mu_x)$ in decreasing order of distance, and stop scanning when, for a point $p_j \in N(\mu_x)$, it is the case that $D(p_j, \mu_x) \leq \frac{1}{2}D(\mu_y, \mu_x)$. By the triangular inequality, any point p_j that satisfies this condition cannot be closer to μ_y than to μ_x . This rule filters out from the scan points whose leader cannot possibly be μ_y , thus significantly speeding up the identification of leaders. Note that all distances between centers in T_i must be available; this implies an added $O(k^2)$ cost for computing and maintaining these distances, which is anyhow dominated by the term $O(nk)$. The gain is that, in practice, fewer than n points need to be scanned at each iteration.

3.1.1 Using medoids

Motivated by the aim of attaining high quality clusters efficiently we have found that we can improve upon a direct application of the M-FPF algorithm by adopting the strategy described below. The M-FPF is applied not to the whole set but only to a random sample of size \sqrt{nk} of the input points (this sample size is suggested in [17]). Afterwards the remaining points are associated to the closest (according to the Generalized Jaccard Distance) center. We obtain improvements in quality

⁸We recall that any vector space is also a metric space, but not vice-versa.

by making an iterative update of the “center” when a new point is associated to a cluster. Within a cluster C_i we find the point a_i furthest from μ_i and the point b_i furthest from a_i (intuitively this is a good approximation to a diametral pair). The medoid m_i is the point in C_i that has the minimum value of the function

$$|D(a_i, x) - D(b_i, x)| + |D(a_i, x) + D(b_i, x)|,$$

over all $x \in C_i$.⁹ When we add a new point p to C_i , we check if the new point should belong to the approximate diametral pair (a_i, b_i) , that is if the distances $D(a_i, p)$ or $D(p, b_i)$ are larger than $D(a_i, b_i)$ and if so we update the approximate diametral pair and m_i accordingly. The association of the remaining points is done with respect to the medoids, rather than the centers. The application of M-PFP on a sample plus the iterative re-computation of medoids gave us a clustering of better quality than simply using M-PFP on the whole input set.

3.2 The candidate words selection algorithm

We select candidate terms for labelling the generated clusters through a modified version of the *information gain* function [4, 37]. Let x be a document taken uniformly at random in the set of documents D . $P(t)$ is the probability that x contains term t , $P(c)$ is the probability that x is in category c . The complementary events are denoted $P(\bar{t}) = 1 - P(t)$ and $P(\bar{c}) = 1 - P(c)$. $P(t, c)$ is the probability that x is in category c and contains term t , $P(\bar{t}, \bar{c})$ is the probability x does not contain t and is not in category c . $P(t, \bar{c})$ is the probability that x contains t but is not in category c . $P(\bar{t}, c)$ is the probability that x does not contain t and is in category c . Clearly being these mutually disjoint events it holds: $P(t, c) + P(\bar{t}, \bar{c}) + P(t, \bar{c}) + P(\bar{t}, c) = 1$. The information gain [4, 37] is the contribution of the four terms:

$$IG(t, c) = \sum_{a \in \{t, \bar{t}\}} \sum_{b \in \{c, \bar{c}\}} P(a, b) \log \frac{P(a, b)}{P(a)P(b)}$$

Intuitively, IG measures the amount of information that each argument contains about the other; when t and c are independent, $IG(t, c) = 0$. This function is often used for feature selection in text classification, where, if $IG(t, c)$ is high, the presence or absence of a term t is deemed to be highly indicative of the membership or non-membership in a category c of the document containing it. In the text classification context, the rationale of including in the sum, aside from the factor that represents the “positive correlation” between the arguments (i.e. the factor $P(t, c) \log \frac{P(t, c)}{P(t)P(c)} + P(\bar{t}, \bar{c}) \log \frac{P(\bar{t}, \bar{c})}{P(\bar{t})P(\bar{c})}$), also the factor that represents their “negative correlation” (i.e. the factor $P(\bar{t}, c) \log \frac{P(\bar{t}, c)}{P(\bar{t})P(c)} + P(t, \bar{c}) \log \frac{P(t, \bar{c})}{P(t)P(\bar{c})}$), is that, if this latter factor has a high value, this means that the absence (resp. presence) of t is highly indicative of the membership (resp. non-membership) of the document in c . That is, the term is useful anyway, although in a “negative” sense.

However, in our context we are interested in terms that *positively describe* the contents of a cluster, and are thus only interested in positive correlation.

$$IG_m(t, c) = P(t, c) \log \frac{P(t, c)}{P(t)P(c)} + P(\bar{t}, \bar{c}) \log \frac{P(\bar{t}, \bar{c})}{P(\bar{t})P(\bar{c})}$$

The keywords selected via IG_m will not be shown directly to the user but will form the input to the label generation phase described in Section 4.

4 Overview of the system

We discuss here in more detail the architecture of Armil. Overall the computation flow is a pipeline consisting in (i) data collection and cleaning, (ii) first-level clus-

⁹This formula mimics in a discrete setting the task of finding the cluster point closest to the median point to the segment (a_i, b_i) in an Euclidean space.

tering, (iii) candidate word extraction for labelling, (iv) second-level clustering, and (v) cluster labelling. Let us review these steps in order.

(1) Querying one or more search engines: The user of *Armil* issues a query string that is re-directed by *Armil* to the selected search engines (at the moment the user can select **Google** and/or **Yahoo!**). As a result *Armil* obtains a list (or several lists) of snippets describing Web pages that the search engines deem relevant to the query. An important system design issue is deciding the type and number of snippet sources to be used as auxiliary search engines. It is well-known that the probability of relevance of a snippet to the user information need quickly decays with the rank of the snippet in the list returned by the search engine. Therefore the need of avoiding low-quality snippets suggests the use of many sources each supplying a low number of high-quality snippets (here high-quality snippet implies both relevance of the snippet to the query and representativeness of the snippet with respect to the associated document). On the other hand, increasing the number of snippet sources raises the pragmatic issue of handling several concurrent threads, the need to detect and handle more duplicates, and the need for a more complex handling of the composite ranking by merging several snippet lists (both globally and within each cluster separately). Since we consider *Armil* a “proof-of-concept” prototype rather than a full-blown service, we have chosen only two (high-quality) sources, **Google** and **Yahoo!**. We limit the total number of snippets to be collected to 200 (of which 80 from **Yahoo!** and 120 from **Google**; these numbers optimize the total waiting time). We produce the initial composite ranking of the merged snippet list by a very simple method, i.e. by alternatively picking snippets from each source list.

(2) Cleaning and filtering: Snippets that are too short or with little informative content (i.e. small number of well formed words) or in non-latin alphabets are filtered out. The input is then filtered by removing non-alphabetic symbols, digits, HTML tags, stop words, and the query terms. These latter are removed since they are likely to be present in every snippet, and thus are going to be useless for the purpose of discriminating different contexts. On the other hand query terms are very significant for the user so they are re-introduced in the label generation phase(7) described below. We then identify the prevalent language of each snippet, which allows us to choose the appropriate stop word list and stemming algorithm. Currently we use the ccTLD (Country Code Top Level Domain) of the url to decide on the prevalent language of a snippet. For the purpose of the experiments we only distinguish between English and Italian. For snippets of English Web pages we use Porter’s stemming algorithm, while for Italian ones we use a simple rule-based stemmer we developed in-house. Currently, no other languages are supported. More advanced language discovery techniques are well-known in literature, and we plan to apply them in the near future, however it should be noticed that, given the on-line requirements, methods too expensive in terms of computational effort should be avoided.

(3) First-level clustering: We build a flat k -clustering representing the first level of the cluster hierarchy, using the M-FPF algorithm with medoids and the Generalized Jaccard Distance as described in Section 3. An important issue is deciding the number k of clusters to create. Currently, this number is fixed to 30, but it is clear that the number of clusters should depend on the query and on the number of snippets found. A general rule is difficult to find, also because the “optimal” number of clusters to be displayed is a function of the goals and preferences of the user. Also, while techniques for automatically determining such optimal number do exist [9], their computational cost is incompatible with the real-time nature of our application. Therefore, besides providing a default value, we allow the user to increase or decrease the value of k to his/her liking. Clusters that contain one snippet only are probably outliers of some sort, and we thus merge them under a single cluster labelled “Other topics”.

(4) Snippets re-ranking: A cluster small enough that the list of its snippets fits in the screen does not require a sophisticated order of presentation. However, in general users are greatly facilitated if the snippets of a cluster are listed in order of their likely importance to the user. Our strategy is to identify a geometric “inner

core” of each cluster and geometric “outliers”. In order to achieve this aim we apply the M-FPF algorithm within each cluster as follows. Since M-FPF is incremental in the parameter k , we increment k up to a value for which it happens that the largest obtained cluster has less than half of the point of the input cluster. This specific group forms the “inner core”, all other points are termed “outliers”. The points in the “core” are shown in the listing before the “outliers”. Within core and outliers points we use a relative ranking obtained by a linear combination of the native ranking generated by the auxiliary search engines. Note that this computation is done only to decide the order of presentation of the snippets at the first level. It should not be confused with the second-level clustering described below.

(5) Candidate words selection: For each cluster we need to determine a set of candidate words for appearing in its label; these will hereafter be referred to as *candidates*. For this purpose, for each word that occurs in the cluster we sum the weights, as defined in Section 3, of all its occurrences in the cluster and pre-select the 10 words with the highest score in each cluster.

We refer to this as *local candidate selection*, since it is done independently for each cluster. For each of the 10 selected terms we compute IG_m , as explained in Section 3.2. The three terms in each cluster with the highest score are chosen as candidates. We refer to this as *global candidate selection*, because the computation of IG_m for a term in a cluster is dependent also on the contents of the other clusters. Global selection has the purpose of obtaining different labels for different clusters. At the end of this procedure, if two clusters have the same signature we merge them, since this is an indication that the target number of clusters k may have been too high for this particular query¹⁰.

(6) Second-level clustering: Although the clustering algorithm could in principle be iterated recursively in each cluster up to an arbitrary number of levels, we limit our algorithm to only two levels, since this is likely to be the maximum depth a user is willing to explore in searching for an answer to his/her information needs¹¹. Second-level clustering is applied to first-level clusters of size larger than a predetermined threshold (at the moment this is fixed to 10 snippets, excluding duplicates). For second-level clustering we adopt a different approach, since metric-based clustering applied at the second level tends to detect a single large “inner core” cluster and several small “outlier” clusters¹². The second-level part of the hierarchy is generated based on the candidate words found for each cluster during the first-level candidate words selection. Calling K the set of three candidate words of a generic cluster, we consider all its subsets as possible signatures for second level clusters. A snippet x is assigned to a signature s if and only if all the signature elements are in x and no candidate in $K \setminus s$ is in x . If a signature is assigned too few snippets (e.g. 1) it is considered as an outlier and it is not shown to the user in the second level. Also, if most of the snippets at the first level end up associated to a single signature, then the second-level clusters are not shown to the user since the second-level subdivision would not be any more useful than the first-level subdivision¹³.

(7) Labelling: Many early clustering systems would use lists of keywords as the output to show to the user. In order to offer a more syntactically pleasant label we decide to give as output well formed phrases or parts of phrases as it is becoming standard in more recent systems (e.g. [7]). We use the candidate keywords just as a basis for generating well-formed phrases that will be shown to the user as real cluster labels¹⁴. Given the title of the Web page contained in the snippet, considered as a sequence of words (this time including stop words) we consider all its contiguous subsequences and we give each subsequence a cumulative score as follows: candidates are given a high positive score (its IG_m score), query words a

¹⁰More precisely, we consider the two original clusters with the same signature as second-level clusters, and we produce for each a different second-level label based on the non-signature keywords of those clusters.

¹¹Vivisimo, for example, uses a two-levels hierarchy

¹²We exploited this phenomenon in the snippet re-ranking step (4).

¹³Since second-level clustering is based on first-level candidate words here we depart slightly from the overall idea of separating clustering and labelling.

¹⁴For example “meaning life”, and “meaning of life” hold the same information but the latter is a grammatically well-formed phrase.

low positive score (set at 0.1), all other words have a negative score (set at -0.2). For labelling a cluster, among all its snippets we select the shortest substring of a snippet title among those having the highest score. This computation can be done efficiently using a dynamic programming approach. The choice of positive and negative weights is to ensure balancing of two conflicting goals: include most of the candidates but few of the non-candidate connecting words in a label.

(8) Exploiting duplicates: Since Armil collects data from several search engines it is possible that the same URL (maybe with a different snippet fragment of text) is present in duplicate. We consider this fact as an indication of the importance of the URL. Therefore, duplicates are accounted for in determining weights and distances. Since clustering is based on title and text, it is possible that the same URL ends up in different clusters, for which it is equally relevant. However, if duplicate snippets appear in the same cluster, they are listed only once with an indication of the two sources. Thus duplicate removal is done just before the presentation to the user.

(9) User Interface: The user interface is important for the success of a Web-based service. We have adopted a scheme common to many search engines and meta-search engines (e.g. Vivisimo), in which the data are shown in ranked list format in the main frame while the list of cluster labels are presented on the left frame as a navigation tree. The interface also allows the user to select the number of clusters, by increasing or decreasing the default value of 30.

5 Anecdotal evidence

In this section we report an analysis of the output of the system on selected queries. In particular we highlight some very high quality clusters obtained. These results were obtained in November 2005.

5.1 Query: armstrong

The query “armstrong” returns 174 snippets organized in 28 clusters, with a clustering time of 0.94 seconds (see Figure 1). Cluster #2, labelled “Armstrong Ceilings”, contains 12 snippets, among which 5 are relative to a company manufacturing ceilings and floors, 4 to companies manufacturing other hardware equipment (e.g. molds, pumps, tools), and one to companies manufacturing medical equipment. Cluster #3, labelled “Lance Armstrong Foundation”, contains 14 snippets, of which 12 are related to the sport champion. Cluster number #4, labelled “Luis Jazz”, contains 20 snippets all related to the well-known jazz musician and mostly in English, while Cluster #6, with the same label, contains 12 snippets of which 9 are relative to the same musician but mostly in Italian, and 3 are relative to other musicians named Armstrong. Cluster #12, labelled “Neil Armstrong”, has 7 snippets of which 6 are related to the well-known astronaut. Cluster #19, labelled “George”, contains 4 snippets, two of which are related to the life of General George Armstrong Custer. Cluster #18, labelled “Edwin Howard Armstrong”, contains 6 snippets, three of which are devoted to Edwin Howard Armstrong, the inventor of the frequency modulation radio broadcasting technique.

5.2 Query: madonna

After adjusting manually the default to 20 clusters, the query “madonna” returns 175 snippets organized in 19 clusters in 0.88 seconds (see Figure 2). Cluster #4, labelled “santuario della madonna del divino”, contains 9 snippets, of which 6 are relative to holy shrines. Cluster #5, labelled “su louise veronica madonna ciccone”, contains 11 snippets, all of them high-quality sites devoted to the well-known popstar. Cluster #7, labelled “Madonna di Campiglio”, contains 13 snippets, of which 11 are related to the Italian ski resort. Cluster #10, labelled “music”, contains 51 snippets, most of which related to the popstar.

```

armstrong (29)
├── armstrong air conditioning and heating (6)
├── ⊕ armstrong ceilings (12)
├── ⊕ lance armstrong foundation (14)
├── ⊕ Louis Jazz (20)
├── steam traps humidification systems (6)
├── ⊕ Louis jazz (12)
├── ⊕ italia (13)
├── salt lake city bed and (5)
├── armstrong county (4)
├── at law with main law (7)
├── guide to the papers of (6)
├── neil armstrong (7)
├── of the village cheese company (4)
├── armstrong biography (9)
├── george (4)
├── armstrong (2)
├── armstrong williams (3)
├── edwin howard armstrong (6)
├── chain market research and consulting (3)
├── sport primi piani ciclismo (7)
├── armstrong harbourtown records (4)
├── armstrong sul sito (4)
├── grande ciclismo tour de france (4)
├── joe armstrong sics (3)
├── armstrong clan society one hundred (1)
├── activity based costing (2)
├── willkommen bei armstrong (2)
├── tim armstrong (2)
└── Other topics (2)

```

Figure 1: Labels generated by Armil for the query “armstrong”.

5.3 Query: allergy

The query “allergy” returns 100 snippets organized in 27 clusters in 0.94 seconds (see Figure 3).

Cluster #3 labelled “Asthma” contains 8 snippets, of which 6 related to asthma and allergy. Cluster #5 labelled “Allergy Journal” contains 4 snippets, all of which relate to periodic medical publications. Cluster #6 labelled “guide to allergies and allergy” contains 7 snippets, of which at least 5 offer general advise on allergy treatments. Cluster #11 labelled “food allergy” contains 5 snippets, all of which related to food allergies.

6 Experimental evaluation of the clustering algorithm

We have performed experiments aimed at assessing the performance of the clustering algorithm underlying Armil.

Assessing the quality of a clustering with respect to a pre-classified collection of items (ground truth) is a well-known problem in IR and objective testing methodologies based on information-theoretic principles have been established. We will compare Armil against the state of the art commercial system Vivisimo and against the classical clustering algorithm k-means (several variants thereof).

6.1 The baseline

As baseline against which to compare the clustering capabilities of Armil, we have chosen Vivisimo¹⁵. Vivisimo is considered an industrial standard in terms of clustering quality and user satisfaction, and in 2001 and 2002 it has won the “best meta-

¹⁵<http://vivisimo.com/>



Figure 2: Labels generated by Armit for the query “madonna”.

search-award” assigned annually by the on-line magazine SearchEngineWatch.com. Vivisimo thus represents a particularly difficult baseline, and it is not known if its clustering quality only depends on an extremely good clustering algorithm, or rather on the use of external knowledge or custom-developed resources. To the best of our knowledge, this is the first published experiment comparing on an objective basis (see below) the clustering quality of an academic prototype and Vivisimo. Vivisimo’s advanced searching feature allows a restriction of the considered auxiliary search engines to a subset of a range of possible auxiliary search engines. For the purpose of our experiment we restrict our source of snippets to the ODP directory.

Visivimo’s clustering search engine has been the focus of research in the field of Human Computer Interface and Information systems (see e.g. [32, 18]), with an emphasis on understanding user behavior by analyzing query logs.

6.2 Measuring clustering quality

Following a consolidated practice, in this paper we measure the effectiveness of a clustering system by the degree to which it is able to “correctly” re-classify a set of pre-classified snippets into exactly the same categories without knowing the original category assignment. In other words, given a set $C = \{c_1, \dots, c_k\}$ of categories, and a set Θ of n snippets pre-classified under C , the “ideal” term clustering algorithm is the one that, when asked to cluster Θ into k groups, produces a grouping $C' = \{c'_1, \dots, c'_k\}$ such that, for each snippet $s_j \in \Theta$, $s_j \in c_i$ if and only if $s_j \in c'_i$. The original labelling is thus viewed as the latent, hidden structure that the clustering system must discover.

The measure we use is *normalized mutual information* (see e.g. [33, page 110][4], and [26]), i.e.

$$NMI(C, C') = \frac{2}{\log |C| |C'|} \sum_{c \in C} \sum_{c' \in C'} P(c, c') \cdot \log \frac{P(c, c')}{P(c) \cdot P(c')},$$

allergy (27)

- aaaaa american academy of allergy (3)
- hon allergy glossary a (2)
- asthma (3)
- allergy the basics (4)
- allergy journal (4)
- guide to allergies and allergy (7)
- allergy discussion group (2)
- information symptoms medication products treatment (5)
- vaccini allergie da acaro vaccino (2)
- food allergy (5)
- allergy and asthma faq (5)
- current allergy and clinical immunology (5)
- allergy clinic diagnosis and treatment (5)
- allergy sinusitis ws hchenor (1)
- institute of allergy and infectious (2)
- hour local news allergy forecast (2)
- allergy hot lists (4)
- latex allergy (3)
- products allergy bedding austin air (4)
- society for allergy and clinical (2)
- allergy magazine (2)
- allergy society of south africa (2)
- agency index of food recalls (3)
- publication allergy (2)
- allergy and immunology internet (2)
- food allergy news (2)
- Other topics (13)

Figure 3: Labels generated by Armil for the query “allergy”.

where $P(c)$ represents the probability that a randomly selected snippet s_j belongs to c , and $P(c, c')$ represents the probability that a randomly selected snippet s_j belongs to both c and c' . The normalization, achieved by the $\frac{2}{\log |C| |C'|}$ factor, is necessary in order to account for the fact that the cardinalities of C and C' are in general different [4].

Higher values of NMI mean better clustering quality. The clustering produced by Vivisimo has partially overlapping clusters (in our experiments Vivisimo assigned roughly 27% of the snippets to more than one cluster), but NMI is designed for non-overlapping clustering. Therefore, in measuring NMI we eliminate from the ground truth, from the clustering produced by Vivisimo, and from that produced by Armil, the snippets that are present in multiple copies.

However, in order to also consider the ability of the two systems to “correctly” duplicate snippets across overlapping clusters, we have also computed the *normalized complementary entropy* [33, page 108] and [4], in which we have changed the normalization factor so as to take overlapping clusters into account. The entropy of a cluster $c'_i \in C'$ is

$$E(c'_i, C) = \sum_{k=1}^{|C|} - \frac{|c'_i \cap c_k|}{|c_k|} \log \frac{|c'_i \cap c_k|}{|c_k|}.$$

The normalized complementary entropy of c'_i is

$$NCE(c'_i, C) = 1 - \frac{E(c'_i, C)}{\log |C|}.$$

NCE ranges in the interval $[0, 1]$, and a greater value implies better quality of c'_i . The complementary normalized entropy of C' is the weighted average of the contributions of the single clusters in C' . Let $n' = \sum_{l \in C'} |c'_l|$ be the sum of the

	Vivisimo	Armil(40)	Armil(30)
<i>NCE</i>	0.667	0.735 (+10.1%)	0.683 (+2.3%)
<i>NMI</i>	0.400	0.442 (+10.5%)	0.406 (+1.5%)

Table 1: Results of the comparative evaluation.

cardinalities of the clusters of C' . Note that when clusters may overlap it holds that $n' \geq n$. Thus

$$NCE(C', C) = \sum_{i \in 1}^{|C'|} \frac{|c'_i|}{n'} NCE(c'_i, C).$$

NCE values reported in Table 1 are thus obtained on the full set of snippets returned by Vivisimo.

6.3 Establishing the ground truth

The ephemeral nature of the Web is amplified by the fact that search engines have at best a partial view of the available pages relevant to a given query. Moreover search engines must produce a ranking of the retrieved relevant pages and display only the pages of highest relevance. Thus establishing a “ground truth” in a context of the full Web is problematic. Following [14], we have made a series of experiments using as input the snippets resulting from queries issued to the Open Directory Project (ODP – see Footnote 6). The ODP is a searchable Web-based directory consisting of a collection of a few million Web pages (as of today, ODP claims to index 5.1M Web pages) pre-classified into more than 590K categories by a group of volunteer human experts. The classification induced by the ODP labelling scheme gives us an objective “ground truth” against which we can compare the clustering quality of Vivisimo and Armil. In ODP, documents are organized according to a hierarchical ontology. For any snippet we obtain a label for its class by considering only the first two levels of the path on the ODP category tree. For example, if a document belongs to class **Games/Puzzles/Anagrams** and another document belongs to class **Games/Puzzles/Crosswords**, we consider both of them to belong to class **Games/Puzzles**. This coarsification is needed in order to balance the number of classes and the number of snippets returned by a query.

Queries are submitted to Vivisimo, asking it to retrieve pages only from ODP. This is done to ensure that Vivisimo and Armil operate on the same set of snippets, hence to ensure full comparability of the results. The resulting set of snippets is parsed and given as input to Armil. Since Vivisimo does not report the ODP category to which a snippet belongs, for each snippet we perform a query to ODP in order to establish its ODP-category.

6.4 Outcome of the comparative experiment

The queries used in this experiment are the last 30 of those reported in Appendix A (the first 5 have been excluded since too few related snippets are present in ODP). On average, ODP returned 41.2 categories for each query. In Table 1 we report the *NMI* and *NCE* values obtained by Vivisimo and Armil on these data. Vivisimo produced by default about 40 clusters; therefore we have run Armil with a target of 40 clusters (thus with a choice close to that of Vivisimo, and to the actual average number of ODP categories per query) and with 30 (this number is the default used in the user evaluation).

The experiments indicate a substantial improvement of about 10% in terms of cluster quality of Armil(40) with respect to Vivisimo.¹⁶ This improvement is an important result since, as noted in 2005 in [7], “[T]he scientific literature offers several solutions to the web-snippet clustering problem, but unfortunately the attainable

¹⁶For the sake of replicating the experiments all the search results have been cached and are available upon request.

performance is far from the one achieved by Vivisimo.” It should be noted moreover that Vivisimo uses a proprietary algorithm, not in the public domain, which might make extensive use of external knowledge. In contrast our algorithm is open and disclosed to the research community.

7 Comparison of Armil and k -means

7.1 The k -means algorithm and its variants

The set of experiments in this Section compare our Variant of the Furthest-Point-First algorithm for k -center clustering with three recently proposed, fast variants of k -means.

The k -means algorithm (see [8, 23]) can be seen as an iterative cluster quality booster. It takes as input a rough k -clustering (or, more precisely, k candidate centroids) and produces as output another k -clustering, hopefully of better quality. It has been shown [31] that by using the sum of squared Euclidean distances as objective function¹⁷, the procedure converges to a local minimum for the objective function within a finite number of iterations.

The main building blocks of k -means are (i) the generation of the initial k candidate centroids, (ii) the main iteration loop, and (iii) the termination condition. In the main iteration loop, given a set of k centroids, each input point is associated to its closest centroid, and the collection of points associated to a centroid is considered as a cluster. For each cluster, a new centroid that is a (weighted) linear combination of the points belonging to the cluster is recomputed, and a new iteration starts¹⁸. Several termination conditions are possible; e.g. the loop can be terminated after a predetermined number of iterations, or when the variation that the centroids have undergone in the last iteration is below a predetermined threshold.

Given the importance of this algorithm, there is a vast literature that discusses its shortcomings and possible improvements to the basic framework. In particular, one well-known such shortcoming is that some clusters may become empty during the computation. To overcome this problem, following [30] we adopt the “ISO-DATA” technique that, when a cluster becomes empty, splits one of the “largest” clusters so as to keep the number of clusters unchanged.

It is well-known, and our experiments confirm this, that the quality of the initialization (i.e. the choice of the initial k centroids) has a deep impact on the resulting accuracy. Several methods for initializing k -means are compared in [29]. As our baselines we have chosen the three such methods that are most amply cited in the literature while being at the same time relatively simple; we have ruled out more advanced and complex initializations since the possible boost in quality would be paid immediately in terms of computational cost, thus bringing about too slow an algorithm for our intended application.

7.2 Algorithms, baselines, and variants

The four algorithms we compare in our experiments are (i) the clustering algorithm in Armil), and (ii) three recently proposed, fast variants of k -means (RC, RP, MQ) [30]. More in detail:

KC This is our k -Center algorithm with triangular inequality filtering, as described in Section 3.

¹⁷More precisely, this corresponds to partitioning S , at every iteration, into non-overlapping clusters C_1, \dots, C_k and determining their centroids $\mu_1, \dots, \mu_k \in V$ so that the sum of the squares of the inter-cluster point-to-center distances

$$\sum_j \sum_{x \in C_j} (D(x, \mu_j))^2$$

is minimized.

¹⁸Note that k -means is defined on vector spaces but not in general on metric spaces, since in metric spaces linear combinations of points may not be points themselves.

- RP This is *k*-means (as described in Section 7.1) with an initialization based on Random Perturbation: for each dimension d_j of the space, the distribution of the projections on d_j of the data points is computed, along with its mean μ_j and its standard deviation σ_j ; the k initial centroids are obtained through k perturbations, driven by the μ_j 's and σ_j 's, of the centroid of all data points [29].
- MQ This is MacQueen's [25] variant of *k*-means: the initial centroids are randomly chosen among the input points, and the remaining points are assigned one at a time to the nearest centroid, and each such assignment causes the immediate recomputation of the centroid involved.
- RC This is again *k*-means where the initial centroids are Randomly Chosen among the input points and the remaining points are assigned to the closest centroid [8].

7.3 Experimental results

Since Web snippet clustering is used as an on-line support to Web browsing, real-time response is a critical parameter. A clustering phase that introduces a delay comparable to the time needed for just downloading the snippets, thus in fact doubling the user waiting time, is not acceptable for most users. For this reason, instead of using a fixed number of iterations as the termination condition of the *k*-means algorithm, we opt for a fixed time deadline (5 seconds, in our experiments) and we halt the clustering algorithms at the end of the first iteration that has run over the deadline.

In Table 2 we report clustering time and output quality of several variants of *k*-center and *k*-means on a sample of 12 queries subdivided, according to the method of [6], in three broad families: *ambiguous queries* (armstrong, jaguar, mandrake, java), *generic queries* (health, language, machine, music, clusters), and *specific queries* (mickey mouse, olympic games, steven spielberg). In Table 2, for each query group we indicate averages of the number n of snippets found, the average number k of ODP classes to which they are associated, and the time τ (in seconds) used to build the vectorial representations. For each query group and each algorithm we indicate the time used for clustering and the quality of the outcome¹⁹. The main conclusion of these experiments is that comparable levels of quality are achieved by Armil and *k*-means, but Armil is 5 to 10 times faster.

Note that in these experiments we ask our system to partition the snippets into exactly $k = 30$ clusters, which is the default of Armil. Notice that this number is reasonably close to the average number of classes found in ODP for the test queries.

8 User evaluation of the cluster labelling algorithm

Assessing “objectively” the quality of a cluster labelling method is a difficult problem, for which no established methodology has gained a wide acceptance. For this reason a user study is the standard testing methodology. We have set up a user evaluation of the cluster labelling component of Armil in order to have an independent and measurable assessment of its performance. We performed the study on 22 volunteer master students, doctoral students and post-docs in computer science at our departments. (the University of Siena, University of Pisa, and IIT-CNR). The volunteers have all a working knowledge of the English language.

The user interface of Armil has been modified so as to show clusters one-by-one and proceed only when the currently shown cluster has been evaluated. The queries are supplied to the evaluators in a round robin fashion from a list of 35 predefined queries (the queries are listed in appendix). For each query the user must first say whether the query is meaningful to him/her; an evaluator is allowed to evaluate only queries meaningful to him/her. For each cluster we propose three

¹⁹For the sake of replicating the experiments all the search results have been cached and are available at <http://pspl.iit.cnr.it/~mcsoft/fpf/fpf.html>

Algorithm	NCE	NMI	Time (sec.)
Query group = "ambiguous \cup generic \cup specific"			
$avg(n) = 157.4$ $avg_{ODP}(k) = 36.9$ $avg(\tau) = 0.729$			
MQ	0.70	0.47	5.59
PTS	0.69	0.45	5.76
RP	0.64	0.41	26.12
KC	0.70	0.47	0.77
Query group = "ambiguous"			
$avg(n) = 165$ $avg_{ODP}(k) = 38.75$ $avg(\tau) = 0.695$			
MQ	0.70	0.45	5.77
PTS	0.67	0.42	6.0
RP	0.63	0.38	33.3
KC	0.69	0.44	0.70
Query group = "specific"			
$avg(n) = 107.3$ $avg_{ODP}(k) = 26$ $avg(\tau) = 0.509$			
MQ	0.77	0.52	5.46
PTS	0.77	0.52	5.32
RP	0.71	0.46	18.47
KC	0.76	0.50	0.71
Query group = "generic"			
$avg(n) = 181.4$ $avg_{ODP}(k) = 42$ $avg(\tau) = 0.889$			
MQ	0.67	0.46	5.51
PTS	0.65	0.44	5.7
RP	0.61	0.40	24.92
KC	0.67	0.47	0.86

Table 2: Time/Quality Results of Web snippet clustering for ODP snippets with Armil and variants of k-means.

questions: (a) Is the label syntactically well-formed?; (b) Can you guess the content of the cluster from the label?; (c) After inspecting the cluster, do you retrospectively consider the cluster as well described by the label? The evaluator must choose one of three possible answers (Yes; Sort-of; No), and his/her answer is automatically recorded in a database. Question (a) is aimed at assessing the gracefulness of the label produced. Question (b) is aimed at assessing the quality of the label as an instrument predictive of the cluster content. Question (c) is aimed at assessing the correspondence of the label with the content of the cluster. Note that the user cannot inspect the content of the cluster before answering (a) and (b).

Selection of the queries. Similarly to [7, 19], we have randomly selected 35 of the most popular queries submitted to Google in 2004 and 2005²⁰; from the selection we have removed queries (such as e.g. "Spongebob", "Hilary Duff") that, referring to someone or something of regional interest only, were unlikely to be meaningful to our evaluators. The selected queries are listed in Appendix A.

Discussion of the results. Each of the 35 queries has been evaluated by two different evaluators, for a total of 70 query evaluations and 1584 cluster evaluations. The results are displayed in the following table:

	Yes	Sort-of	No
(a)	60.5%	25.5%	14.0%
(b)	50.0%	32.0%	18.0%
(c)	47.0%	38.5%	14.5%

Summing the very positive and the mildly positive answers we can conclude that, in this experiment, 86.0% of the labels are syntactically acceptable, 82.0% of the labels are predictive and 85.5% of the clusters are sufficiently well described by their

²⁰<http://www.google.com/press/zeitgeist.html>

label. By checking the percentages of **No** answers, we can notice that sometimes labels considered non-predictive are nonetheless considered well descriptive of the cluster; we interpret this fact as due to the discovery of meanings of the query string previously unknown to the evaluator.

The correlation matrices in Table 3 show more precisely the correlation between syntax, predictivity and representativeness of the labels. Entries in the top part give the percentage over all answers, and entries in the bottom part give percentage over rows.

	b-Yes	b-Sort-of	b-No
a-Yes	42.67%	12.81%	5.11%
a-Sort-of	5.74%	15.27%	4.41%
a-No	1.64%	3.78%	8.52%
a-Yes	70.41%	21.14%	8.43%
a-Sort-of	22.58%	60.04%	17.36%
a-No	11.76%	27.14%	61.08%

	c-Yes	c-Sort-of	c-No
b-Yes	33.52%	12.81%	3.72%
b-Sort-of	11.36%	16.85%	3.66%
b-No	2.14%	8.90%	7.00%
b-Yes	66.96%	25.59%	7.44%
b-Sort-of	35.64%	52.87%	11.48%
b-No	11.88%	49.30%	38.81%

	c-Yes	c-Sort-of	c-No
a-Yes	35.98%	18.93%	5.68%
a-Sort-of	8.64%	12.81%	3.97%
a-No	2.39%	6.81%	4.73%
a-Yes	59.37%	31.25%	9.37%
a-Sort-of	33.99%	50.37%	15.63%
a-No	17.19%	48.86%	33.93%

Table 3: Correlation tables of questions (a) and (b) (top), (b) and (c) (middle), (a) and (c) (bottom).

The data in Table 3 (top) show that there is a strong correlation between syntactic form and predictivity of the labels, as shown by the fact that in a high percentage of cases the same answer was returned to questions (a) and (b).

The middle and bottom part of Table 3 confirms that while for the positive or mildly positive answers (**Yes**, **Sort-of**) there is a strong correlation between the answers returned to the different questions, it is often the case that a label considered not predictive of the content of the cluster can still be found, after inspection of the cluster, to be representative of the content of the cluster.

8.1 Running times

Our system runs on an AMD Athlon (1Ghz Clock) processor with 750Mb RAM and operating system FreeBSD 4.11 - STABLE. The code was developed in Python V. 2.4.1. Excluding the time needed to download the snippets from the auxiliary search engines, the 35 queries have been clustered and labelled in 0.72 seconds on average; the slowest query took 0.92 seconds.

9 Conclusions and future work

Why is Armil not “yet another clustering search engine”? The debate on how to improve the performance of search engines is at the core of the current research in the area of Web studies, and we believe that so far only the surface of the vein has been uncovered. The main philosophy of the system/experiments we have proposed

follows these lines: (i) principled algorithmic choices are made whenever possible; (ii) clustering is clearly decoupled from labelling; (iii) attention is paid to the trade-off between response time and quality while limiting the response time within limits acceptable by the user; (iv) a comparative study of *Armil* and *Vivisimo* has been performed in order to assess the quality of *Armil*'s clustering phase by means of effectiveness measures commonly used in clustering studies; (v) a user study has been set up in order to obtain an indication of user satisfaction with the produced cluster labelling; (vi) no use of external sources of knowledge is made.

Further research is needed in two main areas. First, we plan to assess to what extent a modicum of external knowledge can improve the system's performance without speed penalties. An example of one possible such improvement is the use of extractors of noun phrases in parsing the snippets. Second, it is possible to introduce in the current pipeline (input snippets are clustered, candidates are extracted, labels are generated) of the architecture a feedback loop by considering the extracted candidates/labels as predefined categories, thus examining which snippets in different clusters are closer to the generated labels. Snippets close to the label of cluster C_x but in a different cluster C_y could be shown on the screen as related also to C_x . This would give the benefits of soft clustering without much computational overload. Finally, methods for detecting automatically the desired number of clusters will be studied.

10 Acknowledgments

We wish to thank two anonymous referees, whose comments helped us in improving the style of the presentation, and the participants in the *Armil* user evaluation experiment.

References

- [1] R. A. Baeza-Yates and B. A. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [2] M. S. Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of STOC-02, 34th Annual ACM Symposium on the Theory of Computing*, pages 380–388, Montreal, CA, 2002.
- [3] D. Cheng, R. Kannan, S. Vempala, and G. Wang. On a recursive spectral algorithm for clustering from pairwise similarities. Technical Report MIT-LCS-TR-906, Massachusetts Institute of Technology, Cambridge, US, 2003. <http://cs-www.cs.yale.edu/homes/kannan/Papers/experimentclustering.pdf>.
- [4] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, New York, US, 1991.
- [5] D. R. Cutting, J. O. Pedersen, D. Karger, and J. W. Tukey. Scatter/Gather: A cluster-based approach to browsing large document collections. In *Proceedings of SIGIR-92, 15th ACM International Conference on Research and Development in Information Retrieval*, pages 318–329, Kobenhavn, DK, 1992.
- [6] E. Di Giacomo, W. Didimo, L. Grilli, and G. Liotta. A topology-driven approach to the design of Web meta-search clustering engines. In *Proceedings of SOFSEM-05, 31st Annual Conference on Current Trends in Theory and Practice of Informatics*, pages 106–116, Liptovský Ján, SK, 2005.
- [7] P. Ferragina and A. Gulli. A personalized search engine based on Web-snippet hierarchical clustering. In *Special Interest Tracks and Poster Proceedings of WWW-05, 14th International Conference on the World Wide Web*, pages 801–810, Chiba, JP, 2005.
- [8] E. Forgy. Cluster analysis of multivariate data: Efficiency vs. interpretability of classifications. *Biometrics*, 21(3):768–780, 1965.
- [9] M. C. C. G. W. Milligan. An examination of procedures for determining the number of clusters in a data set. *Psychometrika*, 50(2):159–179, 1985.

- [10] F. Geraci, M. Pellegrini, P. Pisati, and F. Sebastiani. A scalable algorithm for high-quality clustering of Web snippets. In *Proceedings of SAC-06, 21st ACM Symposium on Applied Computing*, pages 1058–1062, Dijon, FR, 2006.
- [11] F. Geraci, M. Pellegrini, F. Sebastiani, and M. Maggini. Cluster generation and cluster labelling for web snippets. In *Proceedings of the 13th Symposium on String Processing and Information Retrieval (SPIRE 2006)*, pages 25–36, Glasgow, UK., October 2006. Volume 4209 in LNCS.
- [12] T. F. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 38(2/3):293–306, 1985.
- [13] T. H. Haveliwala, A. Gionis, and P. Indyk. Scalable techniques for clustering the web. In *WebDB (Informal Proceedings)*, pages 129–134, 2000.
- [14] T. H. Haveliwala, A. Gionis, D. Klein, and P. Indyk. Evaluating strategies for similarity search on the Web. In *Proceedings of WWW-02, 11th International Conference on the World Wide Web*, pages 432–442, Honolulu, US, 2002.
- [15] M. A. Hearst and J. O. Pedersen. Reexamining the cluster hypothesis: Scatter/Gather on retrieval results. In *Proceedings of SIGIR-96, 19th ACM International Conference on Research and Development in Information Retrieval*, pages 76–84, Zürich, CH, 1996.
- [16] D. S. Hochbaum and D. B. Shmoys. A best possible approximation algorithm for the k -center problem. *Mathematics of Operations Research*, 10(2):180–184, 1985.
- [17] P. Indyk. Sublinear time algorithms for metric space problems. In *Proceedings of STOC-99, ACM Symposium on Theory of Computing*, pages 428–434, 1999.
- [18] S. Koshman, A. Spink, and B. Jansen. Using clusters on the vivisimo web search engine. In *Proceedings of HCI International 2005*, Las Vegas, July 22-27 2005. Lawrence Erlbaum Associates, Mahwah, NJ. SBN/ISSN: 0-8058-5807-5, <http://www.hci-international.org/>.
- [19] K. Kumamuru, R. Lotlikar, S. Roy, K. Singal, and R. Krishnapuram. A hierarchical monothetic document clustering algorithm for summarization and browsing search results. In *Proceedings of WWW-04, 13th International Conference on the World Wide Web*, pages 658–665, New York, NY, 2004.
- [20] Y. Kural, S. Robertson, and S. Jones. Deciphering cluster representations. *Information Processing and Management*, 37:593–601, 1993.
- [21] Y. Kural, S. Robertson, and S. Jones. Clustering information retrieval search outputs. In *Proceedings of the 21st BCS IRSG Colloquium on Information Retrieval*, Glasgow, UK, 1999. <http://www.bcs.org/server.php?show=ConWebDoc.4252>.
- [22] D. J. Lawrie and W. B. Croft. Generating hierarchical summaries for Web searches. In *Proceedings of SIGIR-03, 26th ACM International Conference on Research and Development in Information Retrieval*, pages 457–458, 2003.
- [23] S. Lloyd. Least squares quantization in PCM. Technical report, Bell Laboratories, 1957. Reprinted in *IEEE Transactions on Information Theory* IT-28(2), 1982, pp. 129–137.
- [24] Y. Maarek, R. Fagin, I. Ben-Shaul, and D. Pelleg. Ephemeral document clustering for Web applications. Technical Report RJ 10186, IBM, San Jose, US, 2000. <http://citeseer.ist.psu.edu/maarek00ephemeral.html>.
- [25] J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 281–297, 1967.
- [26] M. Meila. Comparing clusterings. In *Proceedings of COLT '03*, volume 2777 of LNCS, pages 173–187, 2003.
- [27] W. Meng, C. T. Yu, and K.-L. Liu. Building efficient and effective metasearch engines. *ACM Computing Surveys*, 34(1):48–89, 2002.
- [28] S. Osinski and D. Weiss. Conceptual clustering using Lingo algorithm: Evaluation on Open Directory Project data. In *Proceedings of IIPWM-04, 5th Conference on Intelligent Information Processing and Web Mining*, pages 369–377, Zakopane, PL, 2004.
- [29] J. M. Peña, J. A. Lozano, and P. Larrañaga. An empirical comparison of four initialization methods for the k -means algorithm. *Pattern Recognition Letters*, 20(10):1027–1040, 1999.

- [30] S. J. Phillips. Acceleration of k -means and related clustering algorithms. In *Proceedings of ALENEX-02, 4th International Workshop on Algorithm Engineering and Experiments*, pages 166–177, San Francisco, US, 2002.
- [31] S. Selim and M. Ismail. K -means type algorithms: A generalized convergence theorem and characterization of local optimality. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6(1):81–87, 1984.
- [32] A. Spink, S. Koshman, B. Jansen, M. Park, and C. Field. Multitasking web search on vivisimo.com. In *Proceedings of IEEE ITCC*, volume 2, pages 486–490, Las Vegas, April 2005.
- [33] A. Strehl. *Relationship-based Clustering and Cluster Ensembles for High-dimensional Data Mining*. PhD thesis, University of Texas, Austin, US, 2002.
- [34] A. Strehl, J. Ghosh, and R. J. Mooney. Impact of similarity measures on Web page clustering. In *Proceedings of the AAAI Workshop on AI for Web Search*, pages 58–64, Austin, US, 2000.
- [35] A. Tombros, R. Villa, and C. J. van Rijsbergen. The effectiveness of query-specific hierarchic clustering in information retrieval. *Information Processing and Management*, 38(4):559–582, 2002.
- [36] G. N. W. Hsu. Easy and hard bottleneck location problems. *Discrete and Applied Mathematics*, 1:209–216, 1979.
- [37] Y. Yang and J. O. Pedersen. A comparative study on feature selection in text categorization. In *Proceedings of ICML-97, 14th International Conference on Machine Learning*, pages 412–420, Nashville, US, 1997.
- [38] O. Zamir and O. Etzioni. Web document clustering: A feasibility demonstration. In *Proceedings of SIGIR-98, 21st ACM International Conference on Research and Development in Information Retrieval*, pages 46–54, Melbourne, AU, 1998.
- [39] O. Zamir, O. Etzioni, O. Madani, and R. M. Karp. Fast and intuitive clustering of Web documents. In *Proceedings of KDD-97, 3rd International Conference on Knowledge Discovery and Data Mining*, pages 287–290, Newport Beach, US, 1997.
- [40] H.-J. Zeng, Q.-C. He, Z. Chen, W.-Y. Ma, and J. Ma. Learning to cluster Web search results. In *Proceedings of SIGIR-04, 27th ACM International Conference on Research and Development in Information Retrieval*, pages 210–217, Sheffield, UK, 2004.
- [41] D. Zhang and Y. Dong. Semantic, hierarchical, online clustering of Web search results. In *Proceedings of APWEB-04, 6th Asia-Pacific Web Conference*, pages 69–78, Hangzhou, CN, 2004.
- [42] T. Zhang, R. Ramakrishnan, and M. Livny. Birch: an efficient data clustering method for very large databases. In *ACM SIGMOD International Conference on Management of Data*, pages 103–114, Montreal, Canada, June 1996.

A Queries used in the user evaluation

skype, winmx, nintendo revolution, pamela anderson, twin towers, wallpaper, firefox, ipod, tsunami, tour de france, weather, matrix, mp3, new orleans, notre dame, games, britney spears, chat, CNN, iraq, james bond, harry potter, simpsons, south park, baseball, ebay, madonna, star wars, tiger, airbus, oscars, london, pink floyd, armstrong, spiderman.