

Web graph similarity for anomaly detection

Panagiotis Papadimitriou · Ali Dasdan ·
Hector Garcia-Molina

Received: 14 December 2009 / Accepted: 27 January 2010 / Published online: 25 February 2010
© The Brazilian Computer Society 2010

Abstract Web graphs are approximate snapshots of the web, created by search engines. They are essential to monitor the evolution of the web and to compute global properties like PageRank values of web pages. Their continuous monitoring requires a notion of graph similarity to help measure the amount and significance of changes in the evolving web. As a result, these measurements provide means to validate how well search engines acquire content from the web. In this paper, we propose five similarity schemes: three of them we adapted from existing graph similarity measures, and two we adapted from well-known document and vector similarity methods (namely, the shingling method and random projection based method). We empirically evaluate and compare all five schemes using a sequence of web graphs from Yahoo!, and study if the schemes can identify anomalies that may occur due to hardware or other problems.

Keywords Anomaly detection · Graph similarity · Locality sensitive hashing

1 Introduction

A search engine has two groups of components: online and offline. The online group processes user queries in real time and returns search results. The offline group collects content from the web and prepares it for the online part.

The offline group essentially consists of three main components: a crawler component to find and download web pages, a web graph component to create graphs out of pages and their links, and an indexer component to index pages using their content.

The quality of the search results produced by the online components depends on the data and structures generated by the offline components: If the crawler has missed a “significant” number of important pages, then the search results are frequently incomplete. If the web graph is not an accurate representation of what was crawled, then the ranking of web pages in search results is negatively affected. Finally, if the index is not consistent with the crawl and the graph, then the search results may be irrelevant to user queries.

In practice, the quality of the offline components can be affected by a variety of what we call *anomalies*. For instance, a web host that is unavailable at crawl time may cause us to miss its content, as well as the content reachable from that host. Since the offline data is massive, spread over 1000s of machines and files, and gathered and manipulated over relatively long periods of time (hours to days), processor, network or software problems can corrupt parts of the data. Similarly, new policies put in place (e.g., to combat web spam or to for grouping virtual hosts into hosts) may have unexpected effects and lead to invalid data.

It is very hard to detect problems with the offline data simply by examining a single snapshot or instance. For example, how can one tell that an important part of the web is missing? Or that IP addresses were not properly grouped into hosts?

Because of the difficulty of identifying anomalies in a single data snapshot, it is more practical to identify anomalies based on “differences” with previous snapshots. The idea is to compute one or more similarity measures between two snapshots (of the crawl, the web graph or the index) and

P. Papadimitriou (✉) · H. Garcia-Molina
Stanford University, Stanford, CA 94305, USA
e-mail: papadimitriou@stanford.edu

H. Garcia-Molina
e-mail: hector@cs.stanford.edu

A. Dasdan
Yahoo! Inc., Sunnyvale, CA 94089, USA
e-mail: ali_dasdan@yahoo.com

anomalies will result in similarities that are too low (or too high). For instance, one similarity metric for the web graph may reflect how many hosts the two graphs share in common. If many hosts are missing from a new web graph, it probably indicates that something went wrong in the crawl or in the graph building. The similarity between graphs can also be used to tune the frequency of crawls. That is, if the similarity between one graph and the next is very high, it may indicate that it is unnecessary to crawl and build the graph so frequently. If the similarity is too low, we may need to crawl more often.

The challenge in this approach to anomaly detection is in developing similarity metrics that (a) can be computed in a reasonable time and in an automated way, and (b) are useful for detecting the types of anomalies experienced in practice. A metric that is too sensitive or sensitive to differences that do not impact the quality of search results will yield too many false positives. Similarly, a metric that is not sensitive enough will yield too many false negatives. There is also the challenge of selecting the proper thresholds that tell us when similarities are “too high” or “too low.”

In this paper, we focus on anomaly detection for the web graph component. In Sect. 3, we provide examples of anomalies that we target. They are similar to real anomalies we have observed at Yahoo!. Anomaly detection in other components is also as important and challenging but not covered here. However, note that crawl anomalies will often also be reflected as web graph anomalies, so our web graph work will also be helpful for dealing with many crawl issues. To the best of our knowledge, this is the first work in the area that addresses the anomaly detection problem for web graph component through the calculation of similarities or differences between consecutive snapshots.

We note that our work is empirical in nature. We take a variety of existing similarity schemes from other areas and apply them to our problem. The domain specificity or high time complexity of these schemes prevent us from applying them directly to web graphs in general and huge web graphs in particular. Hence, we modify and extend the schemes to make them work on the very large web graphs we have in our dataset.

To summarize, the main contributions of this paper are:

- We propose detecting anomalies through their impact on the web graph rather than their causes (Sect. 3).
- We propose using graph similarity for web graphs as a way to detect anomalies (Sect. 4);
- We present five scalable schemes for graph similarity (Sect. 5).
- We propose a locality sensitive hashing (LSH) function for web graphs in the context of our Signature Similarity scheme (Sect. 5.5);
- We provide experimental results using real web graphs from Yahoo!, and anomalies based on actual problems ob-

served at Yahoo!, and compare all the presented schemes and discuss their pros and cons (Sect. 6).

2 Web graphs

Web graphs are useful to compute properties that need a global view. PageRank is one such property but there are usually hundreds of similar properties.

Web graphs from crawlers are at page level, called *page graphs*. By aggregating the vertices on some property, e.g., host, language, or domain, a hierarchy of web graphs can be created. The resulting graphs are often smaller in size, allowing efficient processing. In this paper, we focus on host-level web graphs, called *host graphs*. For some advantages of host graphs, see [23].

A host graph $G = \{V, E\}$ is a directed, weighted graph whose vertices V correspond to web hosts, and whose weighted edges E aggregate the hyperlinks of web pages in these hosts. G has an edge (u, v) if there is at least one hyperlink from a web page of host u to a web page of host v . Each vertex and edge can have one or more weights. For the purpose of this paper, the weight $w(u, v)$ of an edge (u, v) is equal to the number of hyperlinks from pages of u to pages of v .

A *subgraph* G' of a web graph G is a directed, weighted graph whose vertices and edges are subsets of G . We represent subgraph G' as $G' = \{V', E'\}$, where $V' \subseteq V$ and $E' \subseteq E$. We call G' a *vertex-induced* subgraph of G if G' includes all the edges of E whose endpoints are in V' . In the rest of the paper, we refer to the vertex-induced subgraph of vertices V' simply as the subgraph of G with vertices V' .

Besides its topological features, a web graph has called *properties*. These properties can be numerical (scalars or distributions) or categorical (labels or lists). Some vertex properties that we will focus on are PageRank of a vertex as a “quality” score (computed in a host graph [12]), the list of hosts pointing to the vertex (called its *inlinks*), and the list of hosts the vertex is pointing to (called its *outlinks*).

3 Potential anomalies

Since our goal is anomaly detection, we now give examples of the types of anomalies we are interested in detecting. For illustration, we will use the Tiny Web Graph in Fig. 1(a) that contains 8 vertices and 15 edges. A search engine starts crawling from a web page of host A and discovers the rest of the hosts-vertices of the Tiny Web Graph. For each vertex v the search engine keeps three properties: name of the vertex, outlinks and quality $q(v)$, computed from its PageRank. The graph is stored in a column-oriented representation on a 4-machine computer cluster. The files on each machine

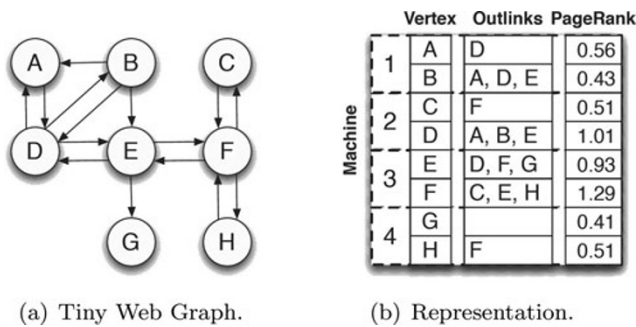


Fig. 1 The Tiny Web Graph and its columnwise representation. The graph on the left has 8 vertices, from A to H, and 15 edges. The representation on the right shows the vertices, their outlinks, their quality or PageRank scores, and how they are stored in a distributed system with 4 machines

correspond to vertex properties: the vertex name, outlinks, and quality. Each machine contains only a fragment (rows) of each file that corresponds to two vertices.

An anomaly occurs when the stored graph representation does not reflect the topology and the properties of the actual web graph at crawl time. An anomaly can be caused by either the problems in the search engine infrastructure, e.g., hardware failures, crawler bugs and web graph data management code bugs, or from problems in the public Internet infrastructure, e.g., host failures and network outages. Although there are ways to handle and fix the problems that arise individually by each of these causes, the problem of monitoring all the infrastructure that is involved in the web graph creation and verifying its proper function is expensive and has no guaranteed overall solution. In this paper we focus on the effects anomalies have on web graphs rather than on the causes of anomalies. Our goal is either to detect an anomaly even if we are unaware of its cause or quantify the effect of an anomaly if we are aware of it. Web graph similarity helps for both goals.

Next we provide a classification of anomalies into three categories based on their impact on the web graph representation with respect to the actual web graph. For each category we provide some anecdotal example causes and illustration in Fig. 2(a).

Missing connected subgraph Anomalies during crawling or machine failures in the storage cluster may result in a

stored web graph that lacks a connected subgraph of the actual web graph. In Fig. 2(a), an anomaly making host F unreachable results in missing the entire component of hosts C, F, and H.

Missing vertices Suppose that machine 4 fails. There are obviously many fault tolerance mechanisms that will recover the data on the failed machine but let us assume the worst case scenario that the data is lost. In this case, we will lose all the information about vertices G and H as shown in Fig. 2(b).

Note that the consequences of a machine failure can be different based on how the vertices are distributed to the machines. If the distribution is such that each machine contains vertices sharing a common property, say, the same country, the failure will lead to the loss of part or all of this country’s data. This case is similar to the removal of a subgraph from the web graph that we discussed above. However, if the distribution is random, the failure will lead to “cuts” all over the web graph. The problem gets more serious if the missing vertices are important vertices in that they have very high rank in quality, or they are major authorities or hubs. For example, failure of machine 3 is more serious than failure of machine 4, since vertices E and F have higher PageRank values than G and H.

Connectivity change Connectivity change refers to anomalies where scalar graph properties such as the total numbers of vertices and edges do not change but for some vertices, the edges these vertices are adjacent to change. This anomaly may be caused by a variety of phenomena. For example, the crawler may misinterpret the hyperlinks of some web pages because of error encoding assumptions. In Fig. 2(c), the crawler failing to discover the edges DA and EG end up inventing two non-existent edges FG and GH.

Note that even if the hyperlinks are correctly collected by the crawler and the corresponding edges are correctly stored, there is no guarantee that the code that is used to access the web graph data is bug-free. For example, modern search engines use column-orientation in flat fragmented files to store graph data, since this approach has many advantages for storing and processing search engine data [3, 7]. However, this approach requires the development of custom code

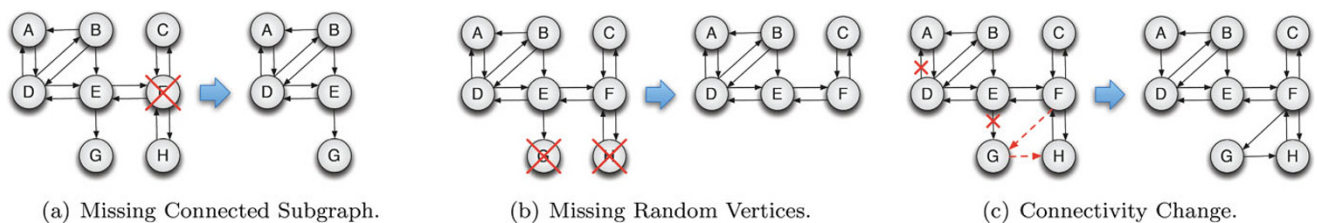


Fig. 2 Examples of anomalies acting on the tiny web graph

to access the graph data, e.g., to fetch vertex names, to join vertices with edges, etc. If there is bug in the code that joins vertices with edges the search engine will rely on a graph with altered edges.

A connectivity change is conceptually simple but is difficult to detect using scalar properties of web graphs. It may also have drastic consequences for a real web graph, since it alters the graph topology that affects the quality scores of the edges in case of PageRank.

4 Web graph similarity

4.1 Problem formulation

We have a sequence of web graphs G_1, \dots, G_n built consecutively and we quantify the changes from one web graph to the next. We do this by computing one or more similarity scores between two consecutive web graphs, G_i and G_{i+1} . Similarity scores, when viewed along the time axis, create a time series. Anomalies can be detected by comparing the similarity score of two graphs against some threshold, or by looking for unusual patterns in the time series.

4.2 Similarity requirements

A similarity function $sim(G, G') \in [0, 1]$ has value 1 if G and G' are identical, and value 0 if G and G' share no common features. The similarity function needs to satisfy the following requirements to be useful for our domain:

1. *Scalability.* The similarity function must produce its result in linear or sublinear time with very small constant factors. This requirement is needed due to the huge size of web graphs.
2. *Sensitivity.* The similarity function must be more sensitive to changes in high-quality vertices and their connectivity. Changes in low-quality vertices are less important and should have a significant impact on the similarity computation only in cases where they affect a large proportion of the web graph vertices.
3. *Coverage.* The similarity function must be sensitive to changes in both topology and the different properties of the web graph. Topological changes include introduction of new hosts, eliminations of existing hosts and changes in the hyperlinks of hosts web pages. As far as the changes in properties, we are mainly concerned about changes in the hosts' quality scores.

The similarity computation schemes we propose are all scalable by our choice and design. We use experimental analysis over different anomaly scenarios to evaluate their sensitivity and coverage for real web graphs.

5 Computing similarities

The problem of comparing graphs or computing their similarity has been an important problem with applications in many areas, from biological networks to web searching. For an overview, see the book [6] and two web documents [22, 24], or search on the Internet using the query “graph similarity”, which returns many useful links in major search engines. Naturally, the diversity of the areas has created different approaches to graph similarity.

After carefully reviewing the families or types of proposed schemes, we identified some families as the most promising for our problem at hand. Within each family, we then developed a particular algorithm that was scalable and would be sensitive to anomalies with high coverage. We discuss these algorithms in Sects. 5.1, 5.2 and 5.3.

In Sects. 5.4 and 5.5, we propose two new approaches to graph similarity [10, 11]. They are inspired by two successful methods for document similarity: the shingling method and the random projection based method. The latter relies on a novel application of a Locality Sensitive Hashing (LSH) function to web graphs.

In the following five subsections, we first describe the family of similarity schemes, then one or more *existing algorithms* in the family, and then our particular *proposed algorithm* adapted from that family. The families of graph similarity schemes that we considered inappropriate for our problem are briefly discussed in Sect. 7.

5.1 Vertex/edge overlap (VEO)

This family uses the rule that “two graphs are similar if they share many vertices and/or edges”. Although this approach seems to ignore the sensitivity and coverage requirements, it is effective in detecting some type of anomalies that we discuss in Sect. 6. Two main ways of computing this kind of similarity is graph edit distance and the Jaccard index [5] (as we apply to graphs). Graph edit distance [6] counts the number of some operations on vertices and edges to transform one graph to the other. The operations consist of insertions, deletions, and in the case of labeled graphs, renamings. The Jaccard index is one way to compute the overlap between two graphs. It is defined as the intersections of vertices/edges divided by the union of vertices/edges.

For this family, we defined the similarity of two graphs G and G' using a form of edit distance as

$$sim_{VEO}(G, G') = 2 \frac{|V \cap V'| + |E \cap E'|}{|V| + |V'| + |E| + |E'|}. \quad (1)$$

We can compute this measure by scanning the vertices (time linear on $|V|$) and checking if each occurs in V' (and similarly for E). If we assume that fetching a specific vertex

or edge by name takes $O(1)$ time, then this similarity computation takes $O(|V| + |V'| + |E| + |E'|)$ time. The algorithm implementing this scheme is called the *Vertex/Edge Overlap Algorithm* (VEO).

5.2 Vertex ranking (VR)

This family uses the rule that “two graphs are similar if the rankings of their vertices are similar”. In this family, the vertices are usually ranked using their qualities, and the similarity of rankings is usually computed using a rank correlation method such as Spearman’s rho (denoted ρ). Although rank correlation is well known in the information retrieval field, its application to graph similarity appears to be new. A very related application, proposed in [4], is to the similarity of ranking algorithms.

Our particular vertex ranking algorithm proceeds as follows. Let $G = (V, E)$ and $G' = (V', E')$ be the two graphs that we want to compare. For each graph we rank the vertices using their quality (scores), producing a sorted list of vertices. Then we find the correlation of the two lists using a modified version of Spearman’s ρ . The modified version satisfies the constraints that (1) the rank correlation needs to be sensitive to quality and (2) it needs to work for partial lists, i.e., the lists that are not permutations of each other. As such, we define ρ as

$$sim_{VR}(G, G') = 1 - \frac{2 \sum_{v \in V \cup V'} w_v \times (\pi_v - \pi'_v)^2}{D}, \quad (2)$$

where π_v and π'_v are the ranks of v in the sorted list for G and G' , respectively, w_v is the quality of v , and D is a normalization factor that limits the maximum value of the fraction to 1.

One subtleness is, what are the quality and rank of a vertex if it does not exist in one of the graphs? Inspired by [13], we resolve it as follows: If a vertex exists in both graphs, its final quality is the average of its qualities from both graphs; otherwise, its final quality is equal to its quality in whichever graph it exists. Moreover, if a vertex does not exist in one graph, say, G' , its rank is equal to $|V'| + 1$. Vertices not in $|G|$ are handled analogously.

Given that vertex quality scores have been precomputed and that the vector of vertices for one graph fits in memory, we can compare two graphs with one scan of their vertex files and an expected running time $O(|V| + |V'|)$. If the second condition fails, the computation can still be done in linear time after the ranked lists are joined on disk based on vertex names. The algorithm implementing this scheme is called the *Vertex Ranking Algorithm* (VR).

5.3 Vertex/edge vector similarity (VS)

This family uses the rule that “two graphs are similar if their node/edge weight vectors are close”. For this family, the weights usually represent quality. More formally, if

v_1, v_2, \dots are the vertices shared by two graphs G and G' , then we build vectors q and q' for these graphs, respectively, where q_i is the quality score of vertex v_i in graph G . Then we compare the two vectors by computing the average difference between q and q' for all i . For some works that use similar approaches, see [6, 20, 21, 25].

For our particular similarity measure, we compare edges, giving each edge a weight that captures the local topology. As a base, for edge (u, v) , we compute weight $\#outlinks(u, v)$ as the number of outlinks from u to v , where these vertices can be pages, hosts, or domains. Then we compute $\gamma(u, v)$ to capture the relative importance of this edge to other edges leaving node u :

$$\gamma(u, v) = \frac{q_u \times \#outlinks(u, v)}{\sum_{\{v':(u,v') \in E\}} \#outlinks(u, v')} \quad (3)$$

where q_u gives the quality of u .

This approach assigns higher weights to edges that start from high quality vertices. Using these γ weights, we calculate the similarity between two graphs G and G' as

$$sim_{VS}(G, G') = 1 - \frac{\sum_{(u,v) \in E \cup E'} \frac{|\gamma(u,v) - \gamma'(u,v)|}{\max(\gamma(u,v), \gamma'(u,v))}}{|E \cup E'|} \quad (4)$$

where the denominator is a normalization factor.

In cases where an edge (u, v) appears in G but not in G' , we let $\gamma'(u, v) = 0$ to set the numerator to one, resulting in a lower similarity between G and G' . The case where (u, v) is in G' but not in G is handled analogously.

This similarity computation runs in $O(|E| + |E'|)$ time. The algorithm implementing this scheme is called the *Vector Similarity Algorithm* (VS).

5.4 Sequence similarity (SeqS)

This family uses the rule that “two graphs are similar if they share many sequences or short paths of vertices and edges”. The algorithms in this family are used to compare objects that are naturally sequenced, e.g., documents that consist of a sequence of words. For example, shingling [5] is frequently used to detect near-duplicate web pages [16].

Because sequence comparison algorithms are efficient and can operate on large inputs, we want to consider them as candidates for our problem. Thus, here we use a sequence comparison scheme, shingling in particular, to compare two graphs. The main challenge for us is converting the graphs into linear sequences that can then be compared using shingling. As far as we know, our proposal is the first application of shingling to graph similarity. However, in a related work [15], shingling was applied to the detection of large dense subgraphs.

We start the description of our Sequence Similarity Algorithm by reviewing the base shingling scheme of [5]. We

then provide a function ϕ that transforms a graph into the input required by shingling, i.e., into a sequence of tokens. The goal for our function ϕ is to produce a sequence that is strongly correlated with the features of the graph so that the sequence similarity reflects graph similarity.

The shingling method works as follows. First, we convert every document into a canonical sequence of tokens $T = \langle t_1, \dots, t_n \rangle$. We then extract all the subsequences of k tokens ($k = 3$ for this work) and fingerprint them using a fingerprinting mechanism that maps each token to a number in a set U . Each of the $n - k + 1$ fingerprints so obtained are called *shingles*. We use $S(T)$ to represent the shingles resulting from the token sequence T . The similarity of two sequences T and T' , and consequently the similarity of the respective documents d and d' , is computed by the ratio $\frac{S(T) \cap S(T')}{S(T) \cup S(T')}$, which is the Jaccard index.

For storage and performance reasons, the shingling scheme actually uses an unbiased estimate to approximate the ratio above. We choose m random permutation functions $\pi_i : U \rightarrow U, 1 \leq i \leq m$ and apply them to the elements of $S(T)$. From the i th permutation $\pi_i(S(T))$ of $S(T)$ we keep the minimum element $\min(\pi_i S(T))$, which is called the i th *minvalue*. In this way, we end up with an m -dimensional vector \mathbf{w} for each document. The similarity estimate for documents d and d' is equal to the percentage of agreeing entries in vectors \mathbf{w} and \mathbf{w}' .

We now turn to the description of our proposed function ϕ that has a graph as its input and produces a sequence elements out of it. In general, the problem of converting a graph into a serial representation such that strongly correlated vertices are placed next to each other is known as the *graph serialization problem* [21]; the exact algorithms for this problem has exponential complexity [21]. To satisfy our similarity requirements from Sect. 4, we propose the following heuristic algorithm called the *walk algorithm* to compute ϕ . This algorithm is similar to the one in [21].

The walk algorithm starts off by visiting the vertex associated with the highest quality, and repeats the following process: (i) Among the unvisited neighbors reachable from the currently visited vertex, visit the neighbor with the highest quality; (ii) If the currently visited vertex does not have any unvisited neighbors, jump to the vertex that has the highest quality among all the unvisited vertices.

To illustrate this algorithm, with Fig. 1, we start with the highest quality vertex F . Then we do a traversal following the links that lead to the highest quality (unvisited) vertices. The highest quality neighbor of F is E , so we visit it next. From there we visit E 's highest quality neighbor, D . After D we visit A . Since A 's neighbors have all been visited, we jump to the highest quality unvisited vertex, C or H . We pick C at random. From C we cannot reach unvisited vertices, so we again jump, this time to H . From H we jump to B and then we jump to the last unvisited node G . The final serialization is F, E, D, A, C, H, B, G .

Given the quality values of the vertices, the walk algorithm runs linearly in the size of the graph. After the walk algorithm, the calculation of shingles and of the vector with the *minvalues* takes time proportional to the number of vertices of the graph. The algorithm implementing this scheme is called the *Sequence Similarity Algorithm* (SeqS).

5.5 Signature similarity (SS)

This family is based on the rule that “two objects are similar if their signatures are similar”. Instead of converting a graph into a sequence, we convert it into a *set* of features, which are then randomly projected into a smaller-dimensionality feature space (signatures) for comparison.

There are many ways to compare sets of features, but here we focus on a scheme called *SimHash*; it is originally developed for high-dimensional vector comparison [8] and is applied to documents comparison [16]. Again, our challenge is in converting our graphs into appropriate sets of features to be input into the SimHash algorithm.

The *SimHash* algorithm works as follows. Initially, a document d is transformed to a set of weighted features $L = \{(t_i, w_i)\}$ where feature t_i is a token of d and w_i is its frequency in d . Tokens are also obtained as in shingling and appear only once in set L . This weighted set can be viewed as a multidimensional vector. The dimensions of the vector space are defined by the distinct features that can be found across all documents. Using this weighted set as input, *SimHash* provides a b -bit fingerprint as follows. Every token t_i is randomly projected into a b -dimensional space by randomly choosing b entries from $\{-w_i, +w_i\}$. We perform this projection using a hash function to obtain the digest for every token t_i . The length of the digest must be equal to b bits (or greater, since we can crop it and reduce it to b bits). We associate $+w_i$ with every 1 and $-w_i$ with every 0 of the digest. The final fingerprint of the document is produced by summing up all the b -dimensional vectors to a single vector h of length b and then converting it to a bit string as follows: the i th entry of h is set to 1 if the entry is positive and 0 if it is negative. The similarity of two documents with fingerprints h and h' is assumed to be equal to the percentage of agreeing bits in h and h' .

The size of the final fingerprint depends on the size of the digest of the hash function we apply to every token. We used 512 bits, which was the largest digest available in the programming language libraries we used (see Sect. 6).

We can view SimHash as a method that provides a similarity estimate between two weighted sets L and L' as

$$\text{sim}_{\text{SimHash}}(L, L') = 1 - \frac{\text{Hamming}(h, h')}{b}, \quad (5)$$

where h and h' are the b -bit vectors corresponding to L and L' , respectively, and $\text{Hamming}(h, h')$ is equal to the

number of bit positions in h and h' for which the corresponding bits are different.

We now move to the problem of defining an appropriate function ϕ that transforms a graph into a weighted set L . Function ϕ returns for each graph a set of weighted features that represent the properties of the graph. Using function ϕ , the similarity between two graphs is then expressed as

$$sim_{SS}(G, G') = sim_{SimHash}(\phi(G), \phi(G')), \tag{6}$$

which is in turn given by (5).

To satisfy our similarity requirements from Sect. 4, we selected the vertices and edges as the features. We assign to a vertex v its quality, and to an edge (u, v) the quality of vertex u normalized by the number of its outlinks, as in Sect. 5.3. For example, the set of weighted features for subgraph G' containing only vertices C, F and H of Fig. 1 is:

$$L(G') = \{(C, 0.51), (CF, 0.51), (F, 1.29), \\ (FC, 1.29 \times 0.5), (FH, 1.29 \times 0.5), \\ (H, 0.51), (HF, 0.51)\}. \tag{7}$$

Since the features of the weighted set are vertices and edges of the graph, we penalize or reward the absence or presence of such common features among different graphs. The use of quality as a weight attributes to every vertex and edge the proper significance from a search engine perspective. The disappearance of a high-quality vertex of a web graph in a new snapshot is penalized more than the disappearance of low-quality vertices.

Assuming again that the quality scores of the vertices of the graph are precomputed, the calculation of its fingerprint requires just one scan of its edge file. This takes time $O(|E|)$. Then, given the fingerprints of two graphs we can compare them in $O(b)$ time, which is $O(1)$ as b is a constant. The algorithm implementing this scheme is called the *Signature Similarity Algorithm* (SS).

6 Experiments and results

We performed experiments to evaluate the similarity measures and algorithms presented in Sect. 5. The experimental section is divided into three subsections.

In Sect. 6.1, we describe the dataset that we used for our experiments and introduce some notation. We also provide some basic statistics about our dataset.

In Sect. 6.2, we experimentally show how similarity varies over time for the different similarity measures and algorithms. This analysis helps us determine the thresholds that identify “significant” changes in a graph.

Finally, in Sect. 6.3, we evaluate how successful the algorithms are in detecting anomalies. In particular, we check

whether significant changes in a web graph affect its similarity score relative to its predecessor graph so that an anomaly can be detected.

6.1 Dataset and setup

For our experiments, we selected a sequence of host web graphs generated by the Yahoo! search engine over a month long period in 2007. These graphs come from crawls that were performed a few times a day to track some “important” web sites with fast changing content such as news sites. Anomaly detection in such graphs is more challenging than in regular web graphs: the fast evolution of the graph, e.g., with addition of new vertices and edges, with elimination of old edges, etc., makes it difficult to determine whether changes in the crawled web graph reflect changes in the actual web graph or are caused by an anomaly.

We refer to a web graph obtained on the X th day of the month as G_X . Since for each day we had multiple graphs, we append lower-case letters to distinguish the ones we chose to use. For example, G_{18a} and G_{18b} are two web graphs built on the 18th day consecutively.

In Fig. 3(a), we show the size of our web graphs on a daily basis. The average size was around several tens of millions. For each day X , we report the number of average number of vertices and edges in graphs G_{Xa}, G_{Xb}, \dots normalized by the number of edges of the first graph of our dataset $V(G_{1a})$.

In Fig. 3(b), we plot the out-degree distribution of vertices of the web graph G_{29b} , normalized by the number of edges in this graph. Both axes are in logarithmic scale. The plot confirms that the number of outlinks per host follows a distribution similar to the power law distribution; a few hosts have a huge number of outlinks whereas almost 60% of hosts have at most one outlink.

In the last figure, Fig. 3(c), we plot the evolution of the web graphs over time. This figure shows the percentage of overlapping vertices and edges between graphs obtained x days apart, as a function of x . Note that x has decimal values due to having multiple graphs per day. We calculated the overlap between the vertices and edges of two graphs G and G' with the Jaccard index.

This figure shows that the overlap for both vertices and edges decreases significantly for graphs less than 5 days apart but far more slowly for larger time differences. The overlap between two graphs built on the same day is higher than 90% but it drops to less than 50% in five days. This rapid decrease in overlap may be attributed to these hosts having fairly fast changing content as in news and blog hosts.

The Yahoo! search engine has reported no anomalies in the software and hardware used to obtain these web graphs. Also, since all of them were successfully used by the search

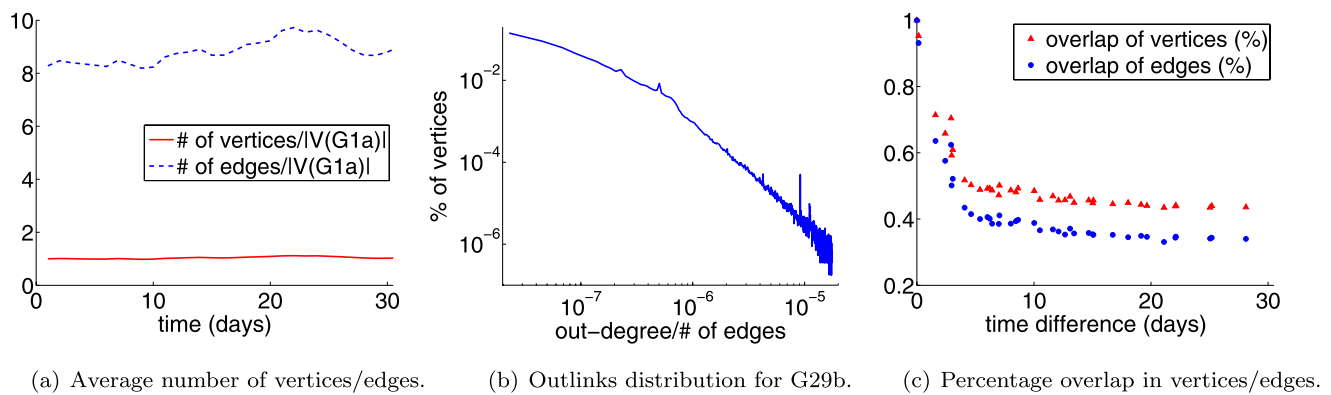


Fig. 3 Statistics of the host web graphs in our dataset. The overlap computation was done using the Jaccard index. The time difference between two graphs refers to the difference between their build times

engine, it is reasonable to assume they had no anomalies. Thus, we may consider the information in Figs. 3(a) and (c) as a reliable evidence of the evolution of the web graphs over time.

We performed the calculation of all the similarity measures locally using a single machine with CloverTownLV 1.86 GHz processor and 6 GB RAM memory. We implemented all the algorithms in Perl and Java. Each similarity computation for one pair of graphs took approximately 20 minutes (wall time), independent of the algorithm used. The running times were dominated by the input–output time spent during the linear scan over the edges file, which is why the running times were similar for all schemes. For this reason, we do not discuss performance further.

6.2 Definition of similarity threshold

A similarity threshold t for an algorithm A is used to identify an anomaly. If two graphs that are expected to be similar, e.g., generated a day apart, have similarity less than t , then we will suspect an anomaly has corrupted one of the graphs. To determine similarity thresholds automatically and in a statistically sound way, we use in our production implementation both non-parametric and parametric methods from time series forecasting, e.g., see [9, 14]. However, the discussion of these methods is beyond the scope of this paper; we refer to [19] for a longer discussion on our selection. Here, we use a simple method with a fixed threshold t for each algorithm that works very well in practice: We set the threshold to the similarity score between two graphs that are built more than one day apart.

6.3 Detecting anomalies

In this section, we try to address the following question: *Can our similarity measures be used to detect common anomalies and quantify their extent (as discussed in Sect. 3)?*

In the following experiments, we select pairs of consecutive web graphs G_{Xa} and G_{Xb} of the same day X and we inject different anomalies into G_{Xb} to get $G_{Xb'}$. We study the impact of the injected anomalies on the similarity between the two web graphs and the sensitivity of the different algorithms to anomalies. Our goal is to see whether we can detect the anomalies in the corrupted $G_{Xb'}$ and estimate their extent, using the similarity score between G_{Xa} and $G_{Xb'}$ and the threshold t from Sect. 6.2. Note that we are comparing $G_{Xb'}$ to G_{Xa} , not to G_{Xb} , since G_{Xb} is the one that has been corrupted and in practice is not available in its uncorrupted form.

The anomalies that we study are missing connected subgraph (Sect. 6.3.1), missing random vertices (Sect. 6.3.2) and connectivity change (Sect. 6.3.3). Our results for each kind of anomaly have two parts:

- We show how each similarity measure captures the effects of an anomaly at different extents. This part of our results provides an insight to the sensitivity and the coverage of each similarity measure that is independent from the application where the measures will be used.
- We also discuss whether a corrupted web graph is acceptable or not and decide whether a similarity measure is good or bad at detecting each anomaly. This part of our results is application dependent and our conclusions are not necessarily valid in different applications.

6.3.1 Missing connected subgraph

We use five pairs of consecutive graphs to simulate the anomaly of missing a connected subgraph: G_{1a} – G_{1b} , G_{18a} – G_{18b} , G_{29a} – G_{29b} , G_{30a} – G_{30b} , and G_{31a} – G_{31b} .

Let a consecutive pair be G_{Xa} – G_{Xb} . We created three corrupted versions of every G_{Xb} by eliminating connected subgraphs that account approximately for 3.5%, 7.0%, 10.5% of its vertices, corresponding to a small, medium, and large amount of corruption, respectively. We refer to the

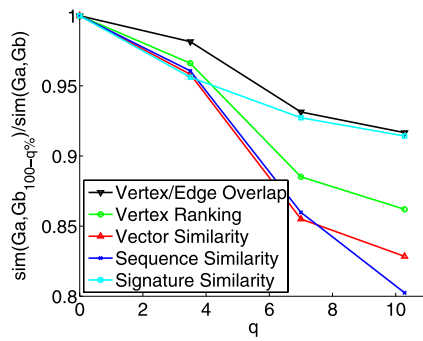


Fig. 4 Similarity of a web graph with its corrupted versions due to missing connected subgraphs. The *x*-axis (log-scale) shows the percentage of the missing vertices and the *y*-axis shows the ratio of the similarity scores

resulting web graphs as $GXB_{100-q\%}$, where q is the percentage of the eliminated vertices of GXB . To eliminate highly connected subgraphs, we removed vertices that belong to specific country domain(s). To create the corrupted version with 3.5%, 7%, 10.5% of missing vertices, we removed vertices of two, one, and three country domains, respectively. We selected domains from the top ranked domains when we sorted them in decreasing order of their contribution to the total number of vertices.

We calculated the similarity of every corrupted version of GXB with the corresponding GXA and compared the result with the similarity of the original web graphs of the pairs, as shown in Fig. 4. Each point on a plot shows the average similarity of an algorithm over the five different pairs. The deviation of similarities around the average is small for all algorithms and we do not report it here to make the plot easier to read.

We observe that the plots of all algorithms decrease monotonically with the extent of changes, but their slopes vary. The Vertex Ranking plot has the steepest descent while the Vertex/Edge Overlap plot decreases more gradually than that of any other algorithm does. Except for the Sequence Similarity and Vertex/Edge Overlap, all the other algorithms successfully detect the two largest corruptions as anomalies. Some algorithms do also detect the smallest corruption as an anomaly but the similarity score is too close to the similarity threshold to make the outcome reliable. Moreover, missing the smallest one is acceptable as the overlap between consecutive web graphs is already in 90–95% range.

6.3.2 Missing random vertices

We used five pairs of consecutive graphs $G1a-G1b$, $G18a-G18b$, $G29a-G29b$, $G30a-G30b$, $G31a-G31b$, and we simulated the removal of random vertices from the GXB graph of each pair. We experimented with different numbers of vertices and vertices of different significance. First, we discuss how we selected the vertices we removed and then we provide the results of our experiments.

We removed from the graph GXB exactly 3.5%, 7.0% and 10.5% of its vertices to create the graphs $GXB_{96.5\%}$, $GXB_{93\%}$ and $GXB_{89.5\%}$, respectively. However, for each percentage, we removed vertices of low, medium, or high quality, where each case is denoted by appending the letters L , M , or H to the name of the resulting graph. For example, we removed three different sets of 7.0% vertices of the web graph $G29b$ to create $G29b_{93\%L}$, $G29b_{93\%M}$ and $G29b_{93\%H}$.

We used *reservoir sampling* to obtain uniformly random samples of fixed size. To bias to the selection of vertices to form $GXB_{q\%}$ with different aggregate quality scores, we performed the following process: we sorted the vertices of GXB in the descending order of their quality scores and divided the resulting list into three subsets. The high, medium, and low quality subsets contained n , $2n$, and $4n$ vertices of GXB so that $n + 2n + 4n = |V_{GXB}|$. The increasing sizes of subsets were motivated by the power law distribution of the quality scores [1]. Finally, we removed vertices with the given percentages from each subset separately.

Before presenting the experimental results, we comment briefly on which corrupted graphs were considered to have significant anomalies. On the one hand, we consider the elimination of high quality vertices as a significant anomaly for all different percentages of the removed vertices (3.5–10.5%). Note that only 3.5% vertices of high quality account for approximately 17% of the aggregate quality score of the web graph. On the other hand, the elimination of medium and low quality vertices should affect the similarity score, but we believe the resulting graph should be considered as not corrupted. Due to the power law distribution of the quality scores, both low and medium quality vertices of any percentage (3.5–10.5%) account for a small portion of the aggregate quality score of the whole graph. In the worst case of removal of 10.5% vertices of medium quality, the removed quality scores account for approximately 6% of the aggregate quality score of the web graph. Given that the overlap of vertices and edges between consecutive graphs with no anomalies is in 90–95% range, e.g., see Fig. 3(c), we can still use the resulting web graph without major problems.

Figure 5 shows our experimental results. We used the same corrupted versions of the web graph GXB to test the different algorithms. Each plot again looks at one similarity algorithm. The *x*-axis of each plot shows the ratio of the aggregate quality score of the corrupted subgraph $GXB_{q\%Y}$, $Y \in \{L, M, H\}$ to the aggregate quality score of GXB . The *y*-axis shows the similarity of a corrupted version of GXB with GXA . There are three plots in each graph, each of them corresponds to the three different percentages q of removed vertices (3.5%, 7.0%, and 10.5%). There are also two additional straight dashed lines that show the similarity score of GXA with the original GXB and the threshold t of each algorithm for comparison purposes.

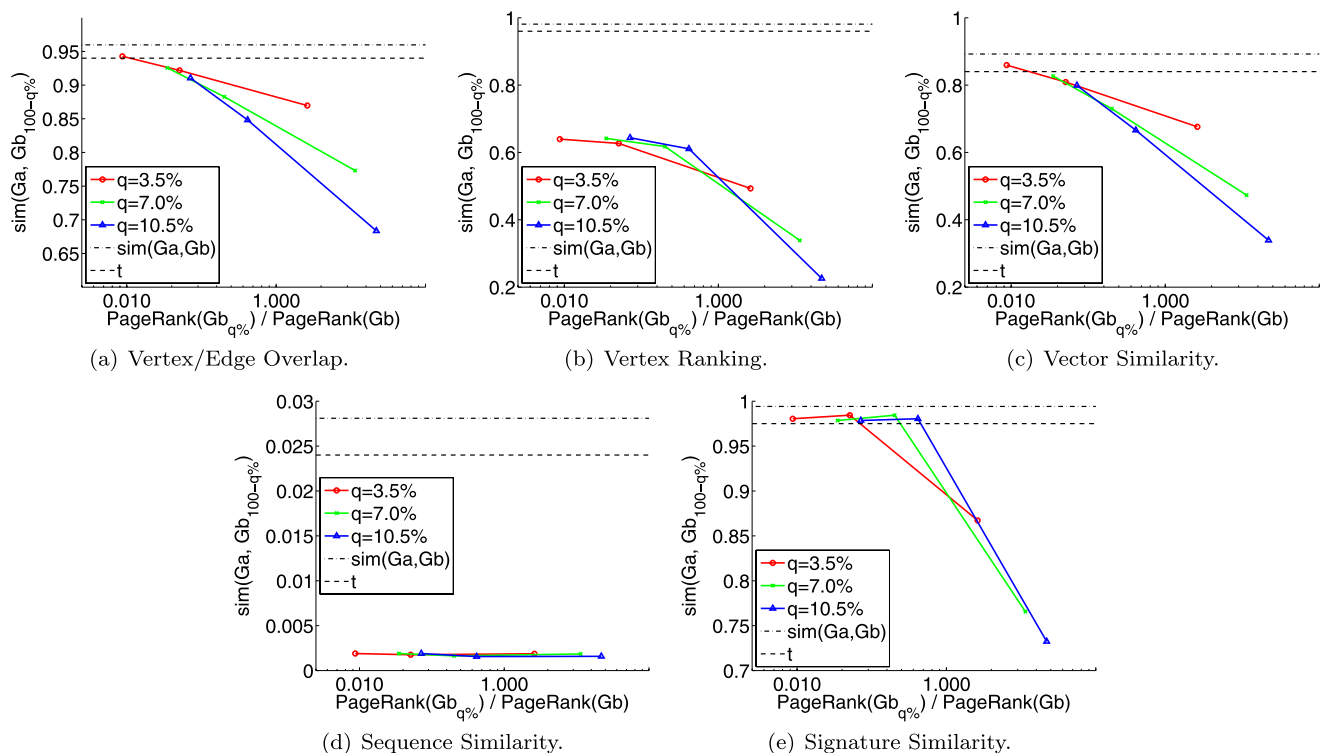


Fig. 5 Similarity of a web graph with its corrupted versions due to missing random vertices. The x -axis shows the percentage of the missing vertices and the y -axis shows the ratio of the similarity scores

We see that Vector Similarity, Vertex/Edge Overlap, and Signature Similarity depend on the significance of changes in the desired way. They give acceptable or nearly acceptable similarity values for changes that affect only insignificant vertices (with Low and Medium quality scores). We see in Fig. 5(c) and in Fig. 5(a) that Vertex Similarity and Vertex/Edge Overlap give similarity scores above or slightly below the threshold for the comparisons of GXa with all corrupted versions of GXb with the missing vertices of Low and Medium quality except from $GXb_{10.5\%M}$. Signature Similarity, shown in Fig. 5(e), gives similarity scores above the threshold for all such comparisons. Although we expected this behavior from Vertex Similarity and Signature Similarity, as they consider explicitly weights for edges and vertices associated with the quality scores of vertices, Vertex/Edge Overlap seems to perform unexpectedly well. We had expected that a similarity measure like Vertex/Edge Overlap that considers only the overlap of features between two graphs would not discriminate between the removal of significant and non-significant vertices. However, the performance of Vertex/Edge Overlap is determined by the way we selected the removed vertices. When we selected vertices with high quality score, we implicitly selected vertices with large number of edges. Since the number of edges is 10:1 to the number of vertices the calculation of the ratio of union to intersection of vertices and edges is dominated by the corresponding ratio of edges. If we had performed the random

selection of the subgraph using random selection of edges, Vertex/Edge Overlap would not have been as successful.

On the contrary, Vertex Ranking and Sequence Similarity showed undesired high sensitivity even in small changes. We see in Fig. 5(b) and Fig. 5(d) that these algorithms give similarity scores below the threshold for all corrupted versions of GXb . These versions include web graphs with only minor changes such as missing 3.5% of low quality vertices of the original graph. This sensitivity would not allow us in a real case to use a web graph where 3.5% of insignificant vertices were missing, even though such a web graph includes most of the desired features to be used by a search engine.

To sum up, Signature Similarity has the desired sensitivity at anomalies of different significance, and that the behavior of Vector Similarity and Vertex/Edge Overlap is close to the desired one. In contrast, Vertex Ranking and Sequence Similarity fail to discern significant from insignificant anomalies.

6.3.3 Connectivity change

We simulated a connectivity change by exchanging outlinks between vertices. In particular we assumed that links are stored by column as discussed in Sect. 3: there is one column for the first outlink of each vertex, another column for the second outlink and so on. To create a corrupted web graph

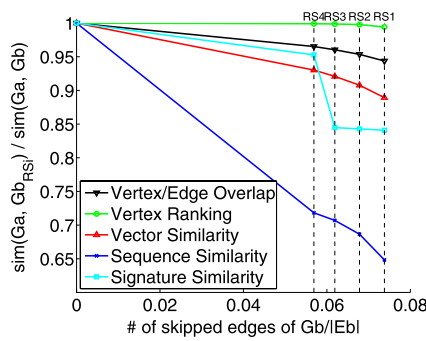


Fig. 6 Similarity of a web graph to its 4 row skipping (RS) versions. The version RS_i refers to the version in which the rows in the i th outlink column were skipped

we skipped the first row of the i th outlinks column, changing the i th outlink of all vertices in the graph. For example, if we skip the 3rd outlinks column in the Tiny Web Graph of Fig. 1(a) we will replace edges DE , EG and FH with EE (self-edge that we can ignore), FG and GH , respectively.

We used five different pairs of consecutive graphs, each pair obtained on a separate day. We simulated the same anomalies in each pair. Let a pair of consecutive subgraphs be $GXA-GXB$. We created four different corrupted versions of GXB by skipping the rows in the i th outlinks column. Let us call such a column a “corrupted column”. We named the resulting corrupted graphs GXB_{RS_i} , $i = 1, \dots, 4$, where i indicates the index of the corrupted column.

Skipping a single row seems like a simple anomaly but it can result in drastic changes in the topology of a web graph. For example, corrupting the first column affects approximately 60% of the hosts, since we saw in Sect. 6.1 that approximately 40% of the hosts have no outlinks. Similarly, corrupting the second column affects approximately 40% of the hosts since 40% of them have no outlinks and 20% of them have exactly one outlink (see Fig. 3(b)). Corrupting the columns with higher indexes affects significantly smaller number of hosts due to the power law distribution of out-degrees. We chose to skip only the first four columns to simulate significant anomalies on the web graphs.

We then calculated the similarities $sim(GXA, GXB_{RS_i})$, $i = 1, \dots, 4$, of the corrupted web graphs to GXA . We present the results only for the pair $G29a-G29b$ since the results of other pairs were similar. Fig. 6 shows the ratios of the computed similarities with the similarity $sim(GXA, GXB)$ of the original web graphs of our pair. The x -axis of Fig. 6 shows the percentage of affected edges of GXB in each case. There are five plots in the plot, one for each algorithm. Every plot has four points that correspond to the similarity of GXA to each of the four GXB_{RS_i} graphs.

These results confirm that Vector Similarity, Sequence Similarity, and Signature Similarity are sensitive to topological changes, as they detect this anomaly. In case of Sequence Similarity, topological changes lead the walk algo-

rithm to different hops while traversing the graph. In cases of the other two algorithms, each edge is assigned weights that are explicitly affected by the significance of its source vertex, which in the case of PageRank as the quality score is also affected by the topology of the graph. Finally, both Vertex/Edge Overlap and Vertex Ranking fail to detect the anomaly, indicating their lack of sensitivity to the topological changes especially when the anomaly results in a small drop in the vertex/edge overlap.

7 Other related work

As discussed in Sect. 5, graph comparison and graph similarity has received a lot of attention in different domains. The book [6] provides a good overview.

Some of the existing similarity schemes were discussed in Sect. 5 as the basis for our own proposed algorithms. Beyond those schemes, there are other schemes that we think are either not applicable to our problem or whose solutions do not scale to the size of the web. Here we summarize two of these well-known approaches.

Graph isomorphism This approach, which is discussed at length in [6], is based on the rule that “two graphs are similar if they are the same (up to labels) or one is contained in the other”. In other words, this approach tries to find out if two graphs are isomorphic to each other (graph isomorphism problem) or one is isomorphic to a subgraph of the other (subgraph isomorphism problem). The former problem is believed to be NP-Complete, and the latter is provably NP-Complete. This approach is not applicable to web graph similarity as isomorphism is too strong a requirement for web graph comparison. Moreover, as the current exact algorithms for this approach all take exponential time, they cannot be applied to real web graphs at all.

Neighborhood similarity This approach is built upon the notion of vertex similarity, which in turn is based on the rule that “two vertices are similar if their neighbors are similar”. This rule appears to have been independently proposed by [2, 17, 18]. These references proposed different ways of computing similarity scores for pairs of vertices. For some of these references, it is also possible to compute a similarity score for graphs once a similarity score for vertex pairs is known. The proposed algorithms have all been iterative, usually involving a fixpoint computation. They are also expensive as their running times are at least quadratic in the number of vertices of one or both graphs. Thus we believe they are not appropriate for our problem either.

Table 1 Effectiveness of similarity schemes in detecting anomalies. Signature Similarity appears to be the most effective among these schemes

Scheme	Anomaly		
	Row skipping	Missing connected	Missing random
Vertex ranking	bad	very good	bad
Vector similarity	good	very good	good
Vertex/edge overlap	bad	very good	good
Sequence similarity	very good	good	bad
Signature similarity	very good	very good	very good

8 Conclusions

Web graphs are at the heart of search engines. Many computers (both internal to the search organization and external to it) are involved in its computation. They are massive and stored on many computers. Because of these complexities many things can go wrong and such anomalies are difficult to detect.

Table 1 summarizes the results of our experiments on real web graphs. Each row of the table refers to a different similarity scheme and each column to a different family of anomalies. Of the five schemes tested, Signature Similarity performed the best followed by Vector Similarity; these two schemes satisfy all three requirements of scalability, sensitivity and coverage. The other three schemes, although they can scale to large web graphs, failed to detect anomalies of certain types.

The common feature of the two successful schemes is the use of appropriately weighted web graph features. The proposed algorithms are independent from the weighting schema that is used and, hence, we believe that they can be effective in anomalies that are not studied here. Our future work will include research on feature selection and weighting schemas for the detection of different types of anomalies in the web graph component of search engine.

References

1. Becchetti L, Castillo C (2006) The distribution of PageRank follows a power-law only for particular values of the damping factor. In: WWW. ACM, New York, pp 941–942
2. Blondel V, Gajardo A, Heymans M, Senellart P, Dooren PV (2004) A measure of similarity between graph vertices: applications to synonym extraction and web searching. *SIAM Rev* 46(4):647–666
3. Boncz P, Kersten M (1999) MIL primitives for querying a fragmented world. *VLDB J* 8(2):101–119
4. Borodin A, Roberts GO, Rosenthal JS, Tsaparas P (2005) Link analysis ranking: algorithms, theory, and experiments. *ACM Trans Internet Technol* 5(1):231–297
5. Broder A, Glassman S, Manasse M, Zweig G (1997) Syntactic clustering of the web. In: WWW, pp 393–404
6. Bunke H, Dickinson PJ, Kraetzl M, Wallis WD (2007) A graph-theoretic approach to enterprise network dynamics. Birkhäuser, Boston
7. Chang F, Dean J, Ghemawat S, Hsieh WC, Wallach DA, Burrows M, Chandra T, Fikes A, Gruber RE (2006) Bigtable: a distributed storage system for structured data. In: OSDI. ACM, New York, pp 205–218
8. Charikar M (2002) Similarity estimation techniques from rounding algorithms. In: STOC. ACM, New York, pp 380–388
9. D’Alberty P, Dasdan A (2009) Non-parametric information-theoretic measures of one-dimensional distribution functions from continuous time series. In: SDM, pp 685–696
10. Dasdan A, Papadimitriou P (2007) Methods and apparatus for computing graph similarity via sequence similarity. US Patent Application, No 20090164411
11. Dasdan A, Papadimitriou P (2007) Methods and apparatus for computing graph similarity via signature similarity, US Patent Application, No 20090150381
12. Eiron N, McCurley KS, Tomlin JA (2004) Ranking the web frontier. In: WWW. ACM, New York, pp 309–318
13. Fagin R, Kumar R, Sivakumar D (2003) Comparing top k lists. *SIAM J Discrete Math* 17(1):134–160
14. Fan J, Yao Q (2005) Nonlinear time series: nonparametric and parametric methods. Springer series in statistics, 2nd edn. Springer, New York
15. Gibson D, Kumar R, Tomkins A (2005) Discovering large dense subgraphs in massive graphs. In: VLDB. ACM, New York, pp 721–732
16. Henzinger MR. Finding near-duplicate web pages: a large-scale evaluation of algorithms. In: SIGIR
17. Jeh G, Widom J (2002) SimRank: a measure of structural-context similarity. In: KDD, Canada
18. Melnik S, Garcia-Molina H, Rahm E (2002) Similarity flooding: a versatile graph matching algorithm and its applications to schema matching. In: ICDE
19. Papadimitriou P, Dasdan A, Garcia-Molina H (2008) Web graph similarity for anomaly detection. Tech Rep 2008-1, Stanford Univ
20. Papadopoulos A, Manolopoulos Y (1999) Structure-based similarity search with graph histograms. In: DEXA/IWOSS. IEEE, New York, pp 174–178
21. Robles-Kelley A, Hancock ER (2005) Graph edit distance from spectral seriation. *IEEE Trans Pattern Anal Mach Intell* 27(3):365–378
22. Seidl T. References for graph similarity. URL: <http://www.dbs.informatik.uni-muenchen.de/~seidl/graphs/>
23. Xue GR, Yang Q, Zeng HJ, Yu Y, Chen Z (2005) Exploiting the hierarchical structure for link analysis. In: SIGIR, New York, NY, USA, pp 186–193
24. Zager L, Verghese G (2005) Graph similarity. URL: <http://lees.mit.edu/lees/presentations/LauraZager.pdf>
25. Zhu P, Wilson RC (2005) A study of graph spectra for comparing graphs. In: MBVC