

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2007-122398

(P2007-122398A)

(43) 公開日 平成19年5月17日(2007.5.17)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30 350C	5B075
G06F 13/00 (2006.01)	G06F 17/30 140	
	G06F 13/00 540F	

審査請求 未請求 請求項の数 5 O L (全 16 頁)

(21) 出願番号	特願2005-313518 (P2005-313518)	(71) 出願人	591089084 中村 健一 東京都足立区谷中5-2-17-207
(22) 出願日	平成17年10月27日(2005.10.27)	(74) 代理人	100120916 弁理士 佐藤 壽見子
		(72) 発明者	中村 健一 東京都足立区谷中5-2-17-207
		Fターム(参考)	5B075 ND03 ND36

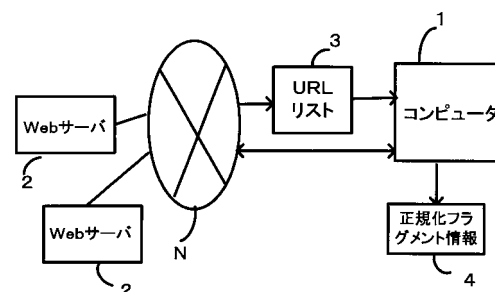
(54) 【発明の名称】 フラグメントの同一性判定方法およびコンピュータプログラム

(57) 【要約】

【課題】 URLが異なるWebページについて、ページ内容の全部又は一部が同一か否かを判定する手段としてフラグメントという概念を導入し、このフラグメントの同一性を判定する方法を提供する。

【解決手段】 内容の異同を判定したいWebページのURLリストを入力情報とし、各URLを指定して該当ページのHTML文書を取得し、各HTML文書から抽出したフラグメントのそれぞれを3つのサブフラグメントに分割し、この3サブフラグメントとURL文字列とを判定要素とすることによって、異なるフラグメント同士との同一性の有無を判定する。判定結果に基づいて、正規化されたフラグメント情報を生成する。この正規化フラグメント情報は、Webページの全部又は一部の内容の異同を推定するために利用することができる。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

インターネットと接続可能なコンピュータが、
インターネット上で公開されている Web ページの任意個数の URL を入力情報とするステップと、

前記の各 URL に対応する Web ページの HTML 文書を取得するステップと、
取得した各 HTML 文書から、リンク要素周辺の文字情報であるフラグメントを抽出するステップと、

抽出した各フラグメントを、リンク要素直前の前フラグメント、リンク要素内部の中フラグメント、およびリンク要素直後の後フラグメントに 3 分割するとともに、これらの分割された 3 種類のサブフラグメントに、該フラグメントを含む HTML 文書に対応する URL の文字列を加えた 4 つを判定要素とするステップと、

或るフラグメントの 4 つの判定要素のそれぞれについて、他のフラグメントの対応する判定要素との類似度を算出するステップと

算出された 4 つの類似度を総合的に評価した結果に基づいて、フラグメント同士の同一性の有無を判定するステップと、

前記同一性判定結果に基づいて、正規化されたフラグメント情報を生成し、出力するステップ、

とからなることを特徴とするフラグメントの同一性判定方法。

【請求項 2】

請求項 1 に記載のフラグメントの同一性判定方法において、
前フラグメントと後フラグメントのそれぞれを構成する要素については、リンク要素からの距離に応じて重み付けを行い、類似度の評価をすることを特徴とするもの。

【請求項 3】

請求項 1 に記載のフラグメントの同一性判定方法において、
3 種類の各サブフラグメントの一致・不一致を評価するためのしきい値の決定に際し、比較される各フラグメントの URL 文字列同士の類似度を利用することを特徴とするもの。

【請求項 4】

請求項 3 に記載のフラグメントの同一性判定方法において、
前フラグメント、中フラグメントおよび後フラグメントの 3 種類のサブフラグメントのうち、少なくとも 2 種類以上が一致すると判定された場合に、同一性のあるフラグメントであると判定することを特徴とするもの。

【請求項 5】

請求項 1 ~ 4 のいずれか 1 に記載の方法を、コンピュータに実行させるためのコンピュータプログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

URL が異なる Web ページについて、ページ内容の全部又は一部が同一か否かを判定するための方法に関する。

【背景技術】**【0002】**

インターネット上には、大量の Web ページが存在し、ある Web ページから他の Web ページへのリンクを張ることが広く行われている。ある Web ページにアクセスした閲覧者は、その Web ページに張られたリンクをたどって別の Web ページへアクセスし、そのアクセス先の Web ページもまた、別の Web ページへのリンクを張っていると、その Web ページへアクセスするというように、次々と Web ページを閲覧していくことができる。

【0003】

また、Webブラウザとサーバとの通信プロトコルであるHTTPには、リファラーと呼ばれる情報が規定されている。このリファラー情報を利用すると、リンク先ページへのリンクを張っているリンク元ページのURL (Uniform Resource Locator) を認識することが可能である。このリファラー情報を利用して、自己をリンク先とするリンク元一覧の表示を行うWebページも出現し、この一覧から適当なURLを指定して、そのリンク元のページを閲覧することもできる。

なお、本出願人が、URL: `http://sourceforge.jp/projects/refererhound/` にて公開しているプログラムも、このリファラー情報を利用している。

【発明の開示】

10

【発明が解決しようとする課題】

【0004】

このようなリンク元一覧を表示することは、たいへん便利なものであり、共通の関心を持つ者同士の交流手段や、広告宣伝手段としてのインターネットの可能性を広げるものである。

しかし、現行のリンク元一覧表示には、いくつかの問題がある。この問題について、図14、図15を参照しながら、説明する。

【0005】

問題点の第一は、リンク元のURLを偽ることが可能であるということである。

例えば、ある邪悪な意図をもったクライアント103aが、スパムを使って宣伝しているようなサイト102aをリンク元としてWebサイト101にアクセスしてきたように装うことが可能である。リンク先サイト101では、このリンク元サイト102aのURLをリンク元一覧に掲載することもありうる(図15のURL104)。その場合、このような事情を知らずに当該サイト101にアクセスしてきたクライアント103bは、リンク元一覧にあるサイト102aへアクセスしてしまう、といったおそれがある。もし、そのサイト102aが、例えば、非合法的なドラッグを販売するようなサイトであったならば、リンク先サイト101の評判まで落としかねない。

20

【0006】

第二の問題点は、通常、一つのWebページは複数のURL、すなわち別名を持つということである。

30

たとえば、Webブラウザ上で、次の5つのURLのどの1つを指定しても、サイト102bの同じWebページが得られるものとする。

`http://www.web_2b.com`

`http://web_2b.com`

`http://www.web_2b.com/index.html`

`http://web_2b.com/index.html`

`http://214.*.*.*(サイト102bのIPアドレス)`

閲覧者が、このリンク元サイト102bの同一ページを介してサイト101にアクセスしてきた場合、上記のリファラー情報を利用すると、リンク先101は、すべて異なるWebページであると認識し、重複してリンク元一覧に載せてしまいかねない(図15のURL105)。そのため、リンク先101にアクセスしてきた閲覧者は、リンク元一覧に掲載されているのは、それぞれ異なるWebページであると勘違いをした結果、同じ内容のWebページに何度もアクセスすることもありうる。

40

また、事実上同一のWebページが重複して掲載されてしまうと、リンク元一覧が膨大となって極めて見にくいWebページとなり、閲覧者の不便を招くことになる。

【0007】

この別名の問題点は、目下、大流行の兆しが見えるWebログ(以下「ブログ」という)において特に重大である。このブログでは、複数の異なるページに同一の記事が含まれ、それらの複数ページは、みな異なるURLをもつからである。

【0008】

50

第1の問題を解決するためには、リンク元102aのWebページに、自己のサイト101に含まれるページへのリンクが存在するか否かを確認すればよい。もし存在しないならば、そのリファラー情報は虚偽であるので、そのリンク元102aのURLは、リンク元一覧に載せないようにすることで問題の解決が図られる。

【0009】

第2の問題、すなわち、同一Webページが複数のURLを持つという問題を解決するためには、リンク元のURLを指定することによりリンク元ページのHTML文書を入力し、このHTML文書の内容から、すでにリンク元一覧に掲載済みのWebページと同一内容のWebページか否かを判定すればよい。そして、URLが異なっても、ページ内容が同一であれば、掲載済みと判定し、リンク元一覧には載せないようにすることで問題が解決できる。

10

【0010】

本発明は、上記の問題解決のためにフラグメントという概念を導入し、このフラグメントの同一性判定結果を、比較されるWebページ同士の全部又は一部の異同の推定に役立てようとするものである。この点では、URL: `http://sourceforge.jp/projects/refererhound/`にて公開しているプログラムも同様である。

しかし、インターネット上には膨大な量のサイトが、相互にリンクを張り巡らしている現状を考えると、このようなフラグメントの異同判定処理に伴う計算量は無視できない。

そのため、本発明は、この計算量を削減し、上記の問題解決手段を実用に耐え得るものとする。このことも目的とする。

20

【0011】

また、計算量を削減しようとする精度が低下しがちである。そのため、異なるフラグメントを同一と判定したり（第一種の過誤）、同一フラグメントを異なるものと判定したり（第二種の過誤）することを極力回避する手段を講じることも目的とする。

なお、本発明は、第二の問題を解決する手法の提供が主眼であるが、判定処理の過程から、第一の問題の解決手段も提供している。

【課題を解決するための手段】

【0012】

このような目的を達成するために、請求項1に記載の発明の方法は、インターネットと接続可能なコンピュータが、インターネット上で公開されているWebページの任意個数のURLを入力情報とするステップと、前記の各URLに対応するWebページのHTML文書を取得するステップと、取得した各HTML文書から、リンク要素周辺の文字情報であるフラグメントを抽出するステップと、抽出した各フラグメントを、リンク要素直前の前フラグメント、リンク要素内部の中フラグメント、およびリンク要素直後の後フラグメントに3分割するとともに、これらの分割された3種類のサブフラグメントに、該フラグメントを含むHTML文書に対応するURLの文字列を加えた4つを判定要素とするステップと、或るフラグメントの4つの判定要素のそれぞれについて、他のフラグメントの対応する判定要素との類似度を算出するステップと算出された4つの類似度を総合的に評価した結果に基づいて、フラグメント同士の同一性の有無を判定するステップと、前記同一性判定結果に基づいて、正規化されたフラグメント情報を生成し、出力するステップ、とからなることを特徴とする。

30

40

【0013】

「リンク元WebページのURLを取得」するために、通信プロトコルHTTPに規定されているリファラー情報を利用する。この情報によって、コンピュータは、リンク元のURLを認識できる。ただし、本発明にとって不可欠であるのは、リファラー情報ではなく、リンク元Webページを探し出す手段である。このような手段のひとつがリファラー情報であるが、他にも、現在広く普及している検索エンジンには、リンク元ページを探し出す機能を有するものもあり、これを利用してよい。

【0014】

50

「HTML文書」とは、Webページを記述したファイルであって、Webブラウザによって読み込まれ、画面にWebページとして表示されるものをいう。このようなファイルは、HTMLで記述されているものが多いので、「HTML文書」と表現する。

「リンク要素」とは、Webページ上にリンクを作り出す要素をいう。代表的なリンク要素としては、HTML文書において、対をなす<A>タグで始まり、タグで終了する文字列があるが、JavaScript（登録商標）を用いてリンクされているものもリンク要素として判別される。

【0015】

「フラグメント」とは、リンク要素周辺の文字情報を単純な文字列として抽出して得た一つのまとまりであって、抽出すべき全体の文字数や、タグの前後の文字数は任意に決めればよい。

10

本発明では、フラグメントを、<A>タグの前か、タグの内側か、タグの後かによって、前フラグメント、中フラグメントおよび後フラグメントの3つのサブフラグメントに分割する。

なお、以下の説明において、サブフラグメントに分割する前のフラグメント、及びサブフラグメントの両者を区別しないで、単に「フラグメント」ということもある。

【0016】

フラグメントは、本発明の基本となる概念である。このフラグメントを利用した判定方法の特徴は、意味を解釈するのではなく、単なる文字列の一致・不一致を判定するという点にある。たとえば、「明るい」と「あかるい」は同一でないと判定する。これは、字面を見るだけであって、文字列の意味まで考慮していないからである。

20

なお、リンク要素は、一つのHTML文書に複数ありうる。或る一つのリンク先に着目し、このリンク先にリンクを張っているリンク元を抽出したい場合は、自己（＝リンク先）のURL文字列を、URLの値として持つリンク要素のみを本発明の処理対象としてもよい。しかし、特定のリンク先に着目することなく、インターネット上で張り巡らされている複数のリンク元・リンク先についてのフラグメント情報を得ようとするのが目的であれば、一つのHTML文書内のすべてのリンク要素を処理対象とすればよい。

【0017】

「正規化されたフラグメント情報」とは、次の2つの条件（条件1）と（条件2）を満たしているフラグメント情報のことである。すなわち、（条件1）は、フラグメントの同一性を判定した結果、互いに一致するフラグメントが重複出現しないこと、（条件2）はこれらの重複のないフラグメントが、それを含むWebページのURLと多対多で対応していることである。

30

【0018】

このように、請求項1に記載の発明によれば、HTML文書の中からフラグメントを取り出し、他のHTML文書中のフラグメントとの異同を比較し、このフラグメントの比較結果によって、Webページの全部又は一部の異同を推定できる。HTML文書全体ではなく、フラグメントの同一性の有無のみを判定すればよいので、計算量が削減される。しかも、フラグメントを3つに分け、比較的字数の少ない文字列同士について別個に類似度の評価をするので、一層計算量が削減される。また、URL文字列同士の類似性も加味することで、判定の精度も向上する。

40

なお、<http://sourceforge.jp/projects/referhound/>において公開されているプログラムも、本発明と同様、リンク要素周辺の文字情報を単純な文字列として抽出し、その同一性の評価によってリンク元の重複を防いでいる。しかし、このプログラムにおいては、リンク要素周辺の文字情報はサブフラグメントに分割されず、一続きのままで同一性を評価している。このため、本発明の方法と同程度の精度で同一性評価を行うには、本発明の方法に比べて140～300%の計算量を要する。

【0019】

また、上記した目的を達成するために、前フラグメントと後フラグメントのそれぞれを構

50

成する要素については、リンク要素からの距離に応じて重み付けを行い、類似度の評価をすることとしてもよい。

これにより、類似度評価の精度が向上する。なぜならば、リンク要素から離れるに従って、それがノイズ成分である可能性が高く、重み付けによって、ノイズ成分を低く評価できるからである。なお、「ノイズ成分」とは、リンク要素と意味的なつながりのない文字列のことをいう。

フラグメントを構成する「要素」とは、文字あるいは部分文字列をいう。

【0020】

さらに、3種類の各サブフラグメントの一致・不一致を評価するためのしきい値の決定に際し、比較される各フラグメントのURL文字列同士の類似度を利用するようにしてもよい。この場合、前フラグメント、中フラグメントおよび後フラグメントの3種類のサブフラグメントのうち、少なくとも2種類以上が一致すると判定された場合に、同一性のあるフラグメントであると判定するようにしてもよい。

10

【0021】

これにより、サブフラグメントだけでなくURL文字列も含む4つの要素についての類似度を総合的に評価するので、フラグメント同士の同一性評価の精度が向上する。

【0022】

上述のフラグメント同士の同一性判定機能を、コンピュータに実現させるためのコンピュータプログラムも本発明である。

【発明の効果】

20

【0023】

アクセス回数の多いWebページには大量のリンク元が存在する。これらのリンク元の重複チェックのためには、リンク元WebページのHTML文書を既知のリンク元HTML文書と比較しなければならない。本発明は、Webページ全体を比較することなく、HTML文書から抽出したフラグメントの同一性を判定することにより、この比較のための計算量を大幅に削減し、重複チェックを実用に耐えるものとすることができる。

一般的に、計算量の削減と精度とはトレードオフの関係にあるが、本発明は、精度の低下を抑制できる。

【発明を実施するための最良の形態】

【0024】

30

(1) 本発明の実施形態であるシステム例の構成

図1に従い、本発明を実施するシステム例を説明する。

コンピュータ1は、インターネットNに接続可能であり、インターネットNに接続するWebサーバ2等の他のコンピュータと情報の送受信を行うものである。

【0025】

コンピュータ1は、インターネットN上に張り巡らされたリンクのリンク元及びリンク先のURLリスト3を入手できる。この入手方法は、リファラー情報の利用によるものでも、現在普及している各種検索エンジンから提供を受けるものでも何でもよい。コンピュータ1は、前記のURLリスト3を入力とし、フラグメントの異同判定処理を行った結果を、正規化フラグメント情報4として出力する。

40

なお、リンク先およびリンク元のWebサイトには、個人が公開するブログなども含まれる。

【0026】

コンピュータ1は、図2に示すように、処理部5、記憶部6、インタフェース部7を備え、他に図示しない入力部や出力部も備えている。

【0027】

記憶部6は、本発明をコンピュータ1に実施させるためのプログラムPROGを記憶したプログラム格納部8、URLリスト格納部9、抽出フラグメントテーブル格納部10、正規化フラグメント情報格納部11を備え、他に各種パラメータなども記憶する。

URLリスト格納部9は、入力されたURLリスト3を記憶しておくものであり、フラグ

50

メント同一性判定処理の間のみ記憶するものであってもよい。

抽出フラグメントテーブル格納部 10 は、HTML 文書から抽出したフラグメントを 3 分割したサブフラグメント及び URL を記憶するものであり、フラグメント同一性判定処理の間のみ記憶するものであってもよい。

正規化フラグメント情報格納部 11 は、出力結果である正規化フラグメント情報 4 を記憶するものである。なお、記憶部 6 は、各種補助記憶装置や ROM, RAM などによって実装されている。

【0028】

処理部 5 は、プログラム格納部 8 からプログラム PROG をメモリ上に読み込み、このプログラム PROG の命令コードを実行する。処理部 5 は、図示しない CPU で実現される。

10

インタフェース部 7 は、外部との信号の送受信を行う。

【0029】

(2) 本発明による判定処理の概要

上記のシステム構成のもと、コンピュータ 1 によって、比較対象となるフラグメントの同一性の有無が、どのように判定されるかについて、図 3 の処理フローに従い、説明する。コンピュータ 1 に 0 個以上のリンク先の URL と、1 個以上のリンク元の URL が列挙された URL リスト 3 が入力される (ステップ S1)。コンピュータ 1 は、URL リスト 3 を参照し 1 個ずつリンク元の URL を取り出し、当該 URL を指定して、Web ページ閲覧要求をし、当該ページの HTML 文書を手に入る (ステップ S2)。

20

【0030】

この HTML 文書からリンク先の URL 文字列を検索する (ステップ S3)。

具体的には、HTML 文書中に、リンク先の URL 文字列が、リンク要素の URL の値として出現するか否かを検索する。

リンク先の URL 文字列が見つからなければ (ステップ S4 で No)、この URL に対応するサイトは虚偽のリンク元であると判断して、この URL 文字列についてのフラグメント抽出処理をスキップする。リンク先とは無関係の Web ページを対象に、フラグメントの正規化を行っても意味がないからである。

【0031】

リンク先の URL 文字列が見つれば (ステップ S4 で Yes)、そのリンク要素周辺の文字情報を取り出す (ステップ S5)。取り出された文字情報が、本発明のフラグメントである。取り出されたフラグメントを 3 種類のサブフラグメントに分割し、これらのサブフラグメントを URL 文字列とともに抽出フラグメントテーブル格納部 10 に登録する (ステップ S6)。図 4 に、このテーブルへの登録例を示す。

30

このようにして、すべてのリンク元 URL に対応する HTML 文書から、リンク先 URL を含むフラグメントをすべて取り出す。

【0032】

次に、抽出された各フラグメントについてステップ S7 から S8 の処理を行う。ある一つのフラグメントに着目し、比較対象となるフラグメントを取り出す (ステップ S7)。

比較対象となるのは、異なる URL に対応した HTML 文書に含まれ、かつ、同一のリンク先を持ち、しかも、未だ同一性判定がなされていないフラグメントである。

40

比較対象たりうるフラグメントが取り出されると、同一性判定処理を行う (ステップ S8)。この処理の内容については、後に詳しく説明する。

【0033】

上記の同一性判定結果に基づき、フラグメント情報の正規化を行う (ステップ S9)。この正規化の処理は、ステップ S9 のように、全フラグメントについての同一性判定を終了してから行ってもよいが、ステップ S6 から S7 のループの中で、逐次行ってもよい。

なお、これを後の処理等で活用するために、正規化フラグメント情報格納部 11 に記憶する。

【0034】

50

ここで、ステップ S 9 のフラグメントの正規化について、具体例を挙げて説明する。

図 5 に示すように、リンク元の URL として、URL_A、URL_B 及び URL_C の 3 つがあり、URL_A に対応する HTML 文書には、リンク先 URL を含むフラグメント A 1、A 2、A 3 の 3 つがあり、URL_B に対応する HTML 文書には、フラグメント B 1、B 2、B 3 の 3 つがあり、URL_C に対応する HTML 文書には、フラグメント C 1、C 2 の 2 つがあるものとする。図 5 の同一性判定結果欄の記号が同じものは、互いに同一性があると判定されたフラグメントである。

この例では、図 6 に示すように、互いに同一性のないフラグメントのグループが 4 つある。

第 1 のグループを、フラグメント A 1 で代表させるならば、フラグメント B 2 は不要の情報となる。同様に、第 2 のグループを、フラグメント A 2 で代表させると、フラグメント B 1 とフラグメント C 2 は不要の情報となり、第 3 のグループを、フラグメント B 3 で代表させると、フラグメント C 1 は不要の情報となる。

【0035】

なお、どのフラグメントで代表させるかについての基準として、例えば、次の (a)、(b) が考えられる。すなわち、(a) そのフラグメントに対応するリンク元からの訪問者数が所定時間 (例えば、過去 24 時間) 内で、一番多いフラグメント、(b) 最後に存在を確認してからの時間がもっとも短いフラグメントである。

(a) の基準で一つに決まれば、そのフラグメントを残すことになるが、もし、決まらなければ、(b) の基準で決めることになる。

【0036】

このような、正規化処理の結果を図 7 に示す。この例からもわかるように、正規化処理の結果は、フラグメントと URL との対応は一般に多対多となる。なお、1 対 1、1 対多、多対 1 対応を含むことは言うまでもない。

図 7 で例示するような多対多の対応を、正規化フラグメント情報格納部 11 に格納する際のデータ構造は、実装レベルの問題であり、特に限定しない。ただし、フラグメントには、図 7 の括弧内に示すように、代表されるフラグメントが含まれていた HTML 文書の URL も併せて登録する。その理由は、フラグメント間の同一性判定には、後述するように、URL 文字列の類似度を組み込んでいるからである。もし、正規化したときに、URL も登録しておかないと、いったん正規化した後に見つかったフラグメントとの間で、URL 文字列を含めた判定ができなくなる。

【0037】

以上の処理の結果、ある Web ページをリンク先とするリンク元のフラグメントの同一性が判定できた。同一性あるフラグメントを有する Web ページ同士は、全部または一部の内容が一致すると推定できる。そのため、正規化フラグメント情報を参照することによって、リンク元一覧には、内容が重複する Web ページの URL を重複掲載しないようにすることができる。図 7 の例では、URL_B と URL_C の Web ページは、URL_A の Web ページと内容的な重複があると推定されるので、URL_A だけをリンク元一覧に掲載すればよい。

なお、Web ページでは、リンク元一覧に、ユーザに示すためのテキストを記述することがある。このテキストは、本発明のフラグメントとは無関係である。フラグメントは、類似性評価用の情報にすぎず、ユーザにどのようなテキストを示すかは、別の問題である。

【0038】

以上、本発明によるフラグメントの同一性判定方法の概要を説明した。

しかし、以上の説明は、あくまで例示である。URL リストを入力とし、正規化フラグメント情報を出力する処理を、少ない計算量で、かつ、必要十分な精度で行うことが本質であって、この処理のための具体的なシステム構成等の実装は限定されるものではない。

以下、フラグメントの抽出と、類似度評価の手法について詳しく説明する。

【0039】

(3) 本発明による判定処理の詳細

10

20

30

40

50

(3-1) フラグメントの抽出

自己のURL文字列を、URLの値としてもつリンク要素の周辺から取り出した文字情報をフラグメントとすることは既に述べたが、具体例をあげて詳細に説明する。

文字情報12として、図8のように、

```
a b c <a href="http://www.xaybzc.co.jp/" target="_blank"> d e f </a> g h i
```

が取り出せたとする。HTMLのタグ部分13を除いた“ a b c d e f g h i ”をフラグメントとするのも一つの方法である。あるいは、タグ部分13の内部の属性値14を取り出し、“ a b c _blank d e f g h i ”をフラグメントとするのもよい。

【0040】

ここで重要なのは、リンク元のWebページ全体(HTML文書全体)の同一性を判別するのではなく、フラグメントの同一性を判別するだけであるという点である。2つのWebページのそれぞれを記述するHTML文書には、同一性あるフラグメントが含まれているならば、これらのWebページの全部あるいは一部が一致すると判断しうる。このように、HTML文書全体を比較することなく、リンク要素周辺から抽出したフラグメント同士の比較だけでよい。これだけでも、Webページの全部又は一部の異同判定に要する計算量は削減されるのであるが、次に示すようなフラグメントの3分割によって、更なる計算量の削減が図られる。

10

【0041】

図9は、文字情報12を3つのサブフラグメントに分割する例を示すものである。文字情報12は、図8と同一である。

20

リンク要素周辺の文字情報12の属性内容と要素内容を、その順番を保持して抽出し、リンク要素の前・中・後を区別して保存する。

文字情報12の開始タグの直前の文字列 a b c を前フラグメント15とし、終了タグの直後の文字列 g h i を後フラグメント17とする。前フラグメント15と後フラグメント17で挟まれた部分

```
<a href="http://www.xaybzc.co.jp/" target="_blank"> d e f </a>
```

からタグの内部の属性値 _blank とタグ外部の文字列 d e f を取り出し、_blank D E F を中フラグメント16として保存する。

【0042】

このようにして、3種類のサブフラグメントが抽出され、これらをフラグメントの同一性判定のための判定要素とする。

30

さらに、リンク元のURL文字列をも判定要素の一つに加える。

なお、4つの判定要素の一つにURL文字列を含めたのは、同じ内容を含むWebページは似通った文字列からなるURLを持つことが多いので、URL文字列の類似度を前・中・後フラグメントの一致・不一致判定のためのしきい値決定に利用するためである。この点については、後に詳しく説明する。

【0043】

(3-2) 文字列同士の類似性

文字列同士を比較する場合、各文字列を構成する個々の文字ごとに比較してもよい。しかし、計算量削減という観点からは、次に述べるように、元の文字列を部分文字列群に分割することが好ましい。

40

また、類似度の算出は、各部分文字列同士を比較して差分を求め、差分の長さで分割前の原文字列の長さとの比に基づいて行うとよい。ここで、「差分の長さ」とは、完全一致しない部分文字列の文字数をいう。「差分の長さ」と原文字列の長さとの比」とは、文字列同士の不一致度であり、この不一致度を1から引いたものが類似度である。

なお、差分を求めるためのアルゴリズムは、すでに公知のものがあるので、これを利用する。〔参考文献：James W. Hunt and Thomas G. Szymanski: A Fast Algorithm for Computing Longest Common Subsequences, Communications of the ACM, vol. 20, no. 5, pp. 350-353, May 1977

50

]

【0044】

(3-2-1) 剰余分割

図10に従い、類似性評価の基本的な手法について説明する。

フラグメント $f r 1$ は、 $C 1$ 、 $C 2$ 、 \dots 、 $C 10$ の10個の文字が順に並んだ文字列であり、フラグメント $f r 2$ は、 $D 1$ 、 $D 2$ 、 \dots 、 $D 11$ の11個の文字が順に並んだ文字列であるとする。

各文字列を、文字コードが5で割り切れる文字のところ分割し、部分文字列のリストに変換する。なお、この分割の仕方を、以下「剰余分割」という。

フラグメント $f r 1$ は、印で囲んだ $C 5$ 、 $C 10$ を剰余が0になる文字とすると、部分文字列 $G 1$ 、 $G 2$ に分割できる。このような剰余分割をフラグメント $f r 2$ についても行うと、部分文字列 $H 1$ 、 $H 2$ と $H 3$ に分割できる。 10

【0045】

異同の判定は、部分文字列のリスト $\{G 1, G 2\}$ と $\{H 1, H 2, H 3\}$ 同士でおこなえばよい。このときの比較演算回数は、両フラグメントの部分文字列の個数の積 ($= 2 \times 3$) である。文字単位で比較する場合は、両フラグメントの文字数の積 ($= 10 \times 11$) だけ、比較演算をしなければならないことを考えると、計算量が大幅に削減される。

【0046】

剰余分割のメリットは、計算量の削減だけではない。剰余分割によれば、文字の挿入・離脱があっても、剰余がゼロになる文字は変わらない、という利点もある。もし、所定個数ずつ原文字列を部分文字列に分割するとしたならば、1文字の挿入あるいは離脱があるだけで、以後の部分文字列がすべて異なってしまう。 20

【0047】

ところで、5で割ることは本質的ではない。この値が大きいほど部分文字列の個数が減り、計算量が削減されるが、類似度評価の精度が低下する。本実施形態では、5という値が適切なトレードオフとなったにすぎない。要するに、自然言語処理のように文字列の持つ意味を考えた、おおがかりな処理をすることなく、機械的に部分文字列に分割すればよいので、計算量が削減できるという点が重要なのである。

【0048】

ただし、部分文字列同士の比較は、1文字が異なるだけで異なる部分文字列と判定されてしまい、文字同士の比較の場合よりも精度が粗くなる。しかし、その代償として、計算量の削減というメリットが得られる。 30

計算量は、小規模な運用では問題にならないが、大規模な運用ではサーバ運用費にかかわる重大な要因となる。

【0049】

(3-2-2) 文字列同士の類似度の算出(基本)

図11に従い、評価の基本的な手法を説明する。

フラグメント $f r 3$ の原文字列 s は、剰余分割により部分文字列 $s 1$ 、 \dots 、 $s 5$ に分割され、フラグメント $f r 4$ の原文字列 t は、剰余分割により部分文字列 $t 1$ 、 \dots 、 $t 6$ に分割されるものとする。 40

文字列 $s t r$ の長さを、 $l e n (s t r)$ と表記し、図11の例では、 $l e n (s) < l e n (t)$ とする。

図11に示す、単に実線で結ばれている $s 1$ と $t 1$ 、 $s 2$ と $t 2$ 、 $s 4$ と $t 5$ および $s 5$ と $t 6$ は互いに同一の部分文字列であり、破線で結ばれている $s 3$ と $t 3$ 、 $s 3$ と $t 4$ は異なる部分文字列であるとする。なお、部分文字列が異なるとは、一文字でも異なる場合をいい、部分文字列の全文字が一致するときは同一とする。

【0050】

上記の例では、フラグメント $f r 3$ の $s 3$ 、フラグメント $f r 4$ の $t 3$ と $t 4$ が差分として取り出される。

両フラグメントの不一致度 D は次の式で得られる。 50

$$D = (\text{len}(t_3) + \text{len}(t_4)) / \text{len}(t)$$

この式からわかるように、フラグメントを構成する文字の個数が異なる場合は、文字数の多いフラグメントを基準に考える。図11の例では、フラグメントfr4が基準となる。なお、類似度(E)は、 $E = 1 - D$ である。

【0051】

(3-2-3) 文字列同士の類似度の算出(重み付けの考慮)

図12に示すように、リンク要素からの距離によって重み付けを行う。ここで、図12は、図11に、重み付け(w_1, \dots, w_6)を追加したにすぎない。

同一フラグメントを構成する文字列であっても、リンク要素に近い部分は同一性成分が並び、リンク要素から遠い部分はノイズ成分が並ぶことが多いと経験的に推測できる。そのため、リンク要素からの距離を考慮して部分文字列に重み付けをする。

図12の例では、フラグメントfr4の部分文字列t1の重み付け係数を w_1 、t2の重み付け係数を w_2 などとする。

この場合の両フラグメントの不一致度Dは次の式で得られる。

$$D = (w_3 * \text{len}(t_3) + w_4 * \text{len}(t_4)) / \text{len}(t)$$

このように重み付けをすることで、計算量削減に伴う精度の低下を抑えることができる。

【0052】

ここで、重み付け係数は1より小さい値で、リンク要素からの距離が大きいほど小さくなる。ノイズ成分を低く評価するためである。

ただし、類似度の算出にあたって、重み付けを考慮するのは、前フラグメントと後フラグメントのみである。中フラグメントは、リンク要素の内側にあり、ノイズ成分の考慮は不要のため、重み付けはしない。

ところで、重み付け係数をいくつにするかは、実装レベルの問題にすぎない。プログラムPROGRAMに定数として記述してもよく、記憶部6に、パラメータとして記憶させてもよい。運用状況によって、適当な値を設定すればよい。

【0053】

なお、上記の式において、分母を $\text{len}(t)$ の代わりに、

$$w_1 * \text{len}(t_1) + w_2 * \text{len}(t_2) + w_3 * \text{len}(t_3) + w_4 * \text{len}(t_4) + w_5 * \text{len}(t_5) + w_6 * \text{len}(t_6)$$

としてもよい。

しかし、実用上、重み付けをしない $\text{len}(t)$ で十分である。

その理由は、たとえリンク要素から遠いところであっても一致しているなら、それはフラグメントが類似していることを強く示唆するからである。

分母を $\text{len}(t)$ とすると、Dの値が小さくなる、すなわち、一致していない場合の不一致評価が軽くなる。これは、類似度Eが大きく評価されることを意味し、同一内容を含むWebページをリンク元として重複掲載したくない、という本発明の目的にかなうものである。

【0054】

(3-3) フラグメントの同一性の判定

図13に従い、3つのサブフラグメントとURL文字列の類似度を総合的に評価してフラグメントの同一性を判定する方法について説明する。

比較対象となる3つのサブフラグメント同士、およびURL文字列同士の4種類の類似度を算出する。

予め、URL文字列同士の類似度に応じて、サブフラグメントの一致・不一致を判定するためのしきい値を定めておく。

図13に示す例では、URL文字列の類似度が0.8を超えるときは、例えば、しきい値を0.7にセットし、0.8以下のときは、例えば、しきい値を0.9にセットする。このように、類似度が高いほど、低いしきい値をセットするのは、URL文字列が似通っているほどフラグメントが一致する可能性が高いからである。

なお、しきい値は、実験段階では、学習データを集めて、適宜設定する。運用段階では、

10

20

30

40

50

日々データが集積していくであろうから、この大量なデータに基づいて、しきい値を更新していくとよい。

【0055】

図13の(ケースA)では、URL文字列の類似度が0.9なので、しきい値は0.7である。前フラグメントと中フラグメントの類似度は、それぞれ0.7を上回っているので、一致しているものとする。したがって、3種類のフラグメントのうち、少なくとも2種類が一致しているので、比較対象のフラグメントは同一性があると判定される。

図13の(ケースB)では、URL文字列の類似度が0.4なので、しきい値は0.9である。しきい値0.9を上回っているのは、中フラグメントだけである。したがって、3種類のサブフラグメントのうち、少なくとも2種類が一致という条件を満たしていないので、比較対象のフラグメントは同一性がないと判定する。

【0056】

このように、フラグメントの同一性判定のために、4つの判定要素の類似度を別個に算定し、これらを総合的に評価した点が、本発明の大きな特徴である。各判定要素は、比較的短い文字列であり、これを部分文字列に分割して差分を求めるので、計算量削減効果がある。

また、URL文字列の類似度をしきい値決定のために利用することは、判定の精度を高めるものである。

【0057】

(4) 計算量削減のための他の手法

同一性あるフラグメントを効率よく見つけるためには、図3のステップS7からS8の処理において、総当り法によることなく、剰余分割された部分文字列を、検索のインデックスとして利用することで、効率よく処理ができる。これについて、以下に説明する。

【0058】

互いに類似性の低いフラグメントの集合があるとき、その集合に含まれる任意のフラグメントは、他のどのフラグメントにもない部分文字列を持つ可能性が高い。こうした部分文字列を「ユニークな部分文字列」と呼ぶことにする。

既知のフラグメントの集合について、ユニークな部分文字列と、その部分文字列を含む既知のフラグメントへのポインタ(上記の実施形態では、抽出フラグメントテーブルのアドレス)を登録した辞書(以下、「指紋辞書」と呼ぶ)をあらかじめ作成しておく。

未知のフラグメントが発見されて、それに類似した既知のフラグメントを探するとき、部分文字列について指紋辞書をあたることで、類似性が高いと見込まれるフラグメントを、抽出フラグメントテーブルから迅速に検索することができる。

【0059】

因みに、既存の自然言語処理では意味的な観点からの利用を前提とするため、インデックスを作成するための分割処理に、形態素解析などの計算量の多い処理を用いている。しかし、本発明では、文字情報の持つ意味的な類似性はまったく考慮しないので、機械的な剰余分割によって得られた部分文字列が、直ちにテーブル検索のためのインデックスとして利用できる。また、同様の機械的な分割であるNグラムに比べても、少ない計算量で目的を達せられる。

【0060】

(5) その他

本発明の同一性判定プログラムは、どのようなコンピュータ言語で記述されていてもよい。Webページを記述したHTML文書中で、当該プログラムが引用されるようになっていてもよい。

【0061】

また、本発明の方法を実施するコンピュータは、1又は複数のリンク先サイトを管理するWebサーバであってもよいが、サイトとは無関係に、フラグメントの正規化処理サービスを行うものであってもよい。

【0062】

10

20

30

40

50

さらに、入力URLリストは、特にリンク先、リンク元という別は必須ではない。リンク先、リンク元を区別することなく、インターネット上に張り巡らされたサイト群から抽出したフラグメントの正規化処理を行う場合もあるからである。その場合、図3のステップS3とS4の処理は不要なことは言うまでもない。

【0063】

コンピュータ1へのURLリスト3の入力の仕方は、インターネットNを介して受信しても、管理者等によるキーボード入力等でも、URL情報を格納した記憶媒体から読み込むものであっても、何でもよい。

【0064】

出力される正規化フラグメント情報4は、コンピュータ1の補助記憶装置に格納するものに限らず、コンピュータ1とは別のデータベースサーバなどに格納してもよい。あるいは、本発明のコンピュータが、他のコンピュータ(PDAや携帯電話も含む)からの依頼により、フラグメント同一性判定サービスを行うならば、自分に接続している記憶装置に格納しなくてもよい。代わりに、サービス依頼元のコンピュータに、正規化フラグメント情報4を送信したり、プリンタ出力したりしてもよい。

10

【0065】

なお、本発明の、フラグメントという概念を利用した同一性判定の手法は、リンク元一覧表示という利用の仕方限定されない。創意工夫次第で、インターネット上のさまざまな場面で活用される可能性を秘めている。

【0066】

要は、以上のように開示された実施の形態はすべての点で例示であって、制限的なものではないということである。したがって、種々の変形が可能である。しかし、その変形が特許請求の範囲に記載された技術思想に基づくものである限り、その変形は本発明の技術的範囲に含まれる。

20

【図面の簡単な説明】

【0067】

【図1】本発明の実施形態のシステム構成例を示す図である。

【図2】本発明の実施形態で用いるコンピュータのブロック図である。

【図3】本発明の実施形態の処理概要を示す流れ図である。

30

【図4】抽出されたフラグメントを分割して格納したテーブル構造を例示する図である。

【図5】フラグメントの正規化を説明するための図である。

【図6】フラグメントの正規化を説明するための図である。

【図7】フラグメントの正規化を説明するための図である。

【図8】フラグメントの抽出を説明するための図である。

【図9】3種類のサブフラグメントの抽出を説明するための図である。

【図10】文字列の剰余分割を説明するための図である。

【図11】文字列を部分文字列のリストに変換した後、差分を抽出することを説明するための図である。

【図12】文字列を部分文字列のリストに変換した後、重み付けを考慮して差分を抽出することを説明するための図である。

40

【図13】フラグメントの異同の判定の仕方を説明するための図である。

【図14】従来技術の説明のためのシステム構成例を示す図である。

【図15】従来技術の説明のためのリンク元一覧表示例を示す図である。 _

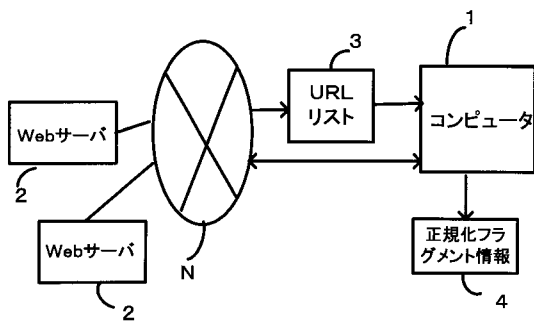
【符号の説明】

【0068】

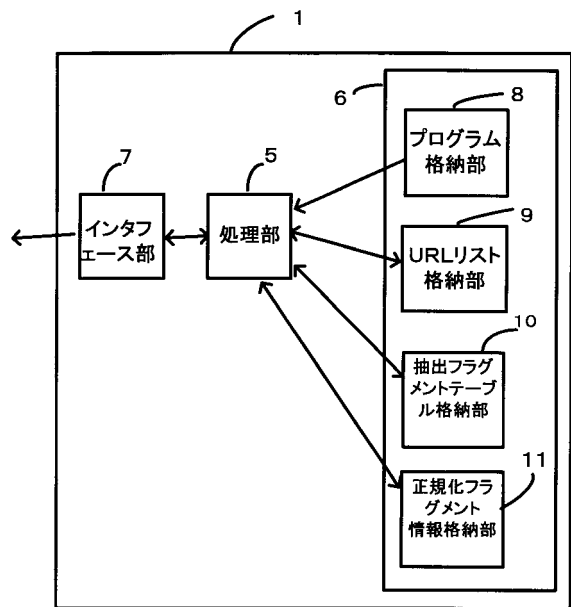
- 1 コンピュータ
- 3 URLリスト
- 4 正規化フラグメント情報
- N インターネット

50

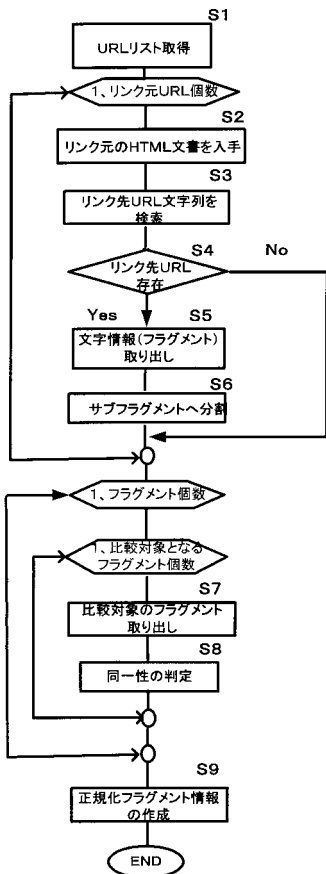
【 図 1 】



【 図 2 】



【 図 3 】



【 図 4 】

レコードNo	前フラグメント	中フラグメント	後フラグメント	URL文字列
1	abc	_blankDEF	ghi	Web_2b.com

【 図 5 】

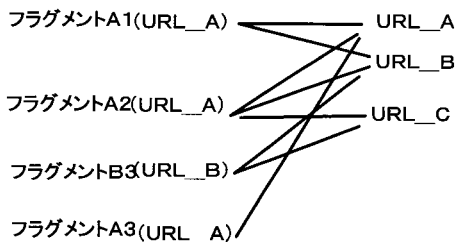
WebページのURL	HTML文書に含まれるフラグメント	同一性判定結果 【注】
URL_A	フラグメントA1	○
	フラグメントA2	△
	フラグメントA3	—
URL_B	フラグメントB1	△
	フラグメントB2	○
	フラグメントB3	◇
URL_C	フラグメントC1	◇
	フラグメントC2	△

【注】 同じ符号は、同一性があると判定されたフラグメントであることを意味する。

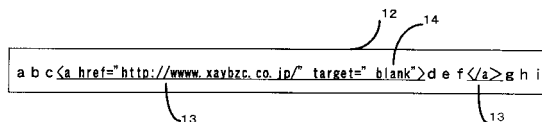
【 図 6 】

グループ No.	同一性判定結果	フラグメント	URL
1	○	フラグメントA1	URL_A
		フラグメントB2	URL_B
2	△	フラグメントA2	URL_A
		フラグメントB1	URL_B
		フラグメントC2	URL_C
3	◇	フラグメントB3	URL_B
		フラグメントC1	URL_C
4	—	フラグメントA3	URL_A

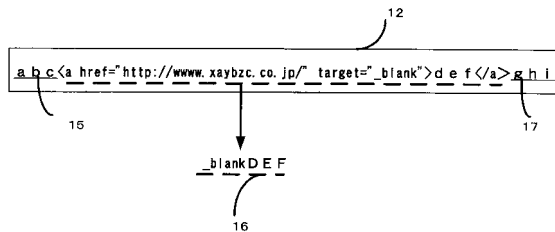
【 図 7 】



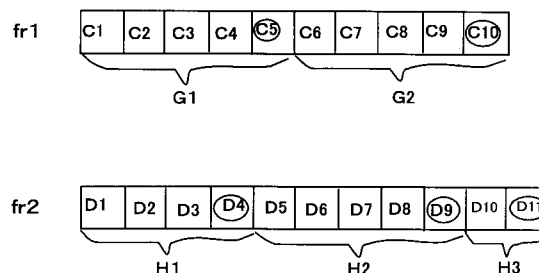
【 図 8 】



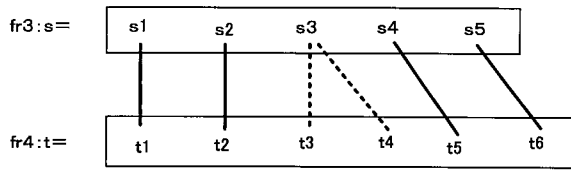
【 図 9 】



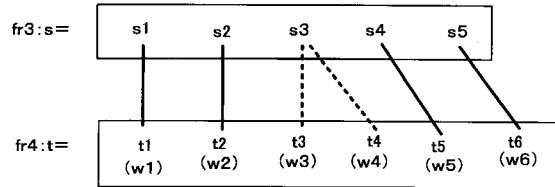
【 図 10 】



【 図 1 1 】



【 図 1 2 】



【 図 1 3 】

(類似度のしきい値)
 URL文字列の類似度 >0.8 ならば、0.7
 URL文字列の類似度 <=0.8 ならば、0.9

(ケースA)

判定要素	類似度	一致(O)、不一致(x)
URL文字列	0.9	しきい値は、0.7
前フラグメント	0.75	O
中フラグメント	0.8	O
後フラグメント	0.5	x

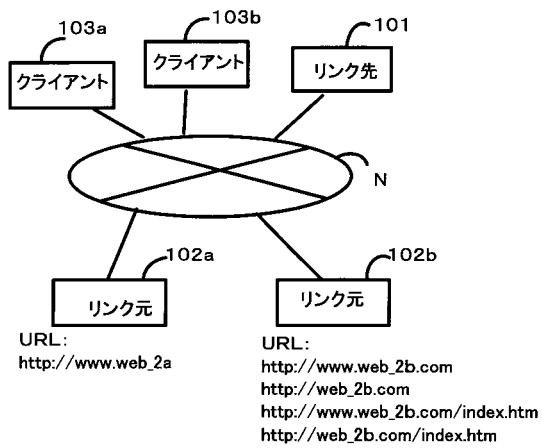
⇒ 判定結果: 同一性あるフラグメント

(ケースB)

判定要素	類似度	一致(O)、不一致(x)
URL文字列	0.4	しきい値は、0.9
前フラグメント	0.8	x
中フラグメント	0.95	O
後フラグメント	0.7	x

⇒ 判定結果: 同一性のないフラグメント

【 図 1 4 】



【 図 1 5 】

