

(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2011-248762

(P2011-248762A)

(43) 公開日 平成23年12月8日(2011.12.8)

(51) Int.Cl.	F I	テーマコード (参考)
G 0 6 F 17/30 (2006.01)	G 0 6 F 17/30 2 1 0 D	5 B 0 7 5
	G 0 6 F 17/30 3 3 0 Z	

審査請求 未請求 請求項の数 8 O L (全 13 頁)

(21) 出願番号	特願2010-123275 (P2010-123275)	(71) 出願人	392026693
(22) 出願日	平成22年5月28日 (2010. 5. 28)		株式会社エヌ・ティ・ティ・ドコモ
			東京都千代田区永田町二丁目 1 1 番 1 号
		(74) 代理人	110000752
			特許業務法人朝日特許事務所
		(72) 発明者	佐々木 純
			東京都千代田区永田町二丁目 1 1 番 1 号
			株式会社エヌ・ティ・ティ・ドコモ内
		F ターム (参考)	5B075 NK02 NR12 PP22 PR03

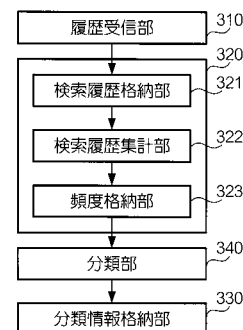
(54) 【発明の名称】 分類装置、コンテンツ検索システム、コンテンツ分類方法、コンテンツ検索方法及びプログラム

(57) 【要約】

【課題】コンテンツを、当該コンテンツの内容によらないで分類できるようにする。

【解決手段】分類情報格納部 3 3 0 は、カテゴリが既知のコンテンツに対応するカテゴリの情報を格納する。履歴受信部 3 1 0 は、コンテンツの検索履歴を受信する。検索履歴は、コンテンツを検索するために用いられた検索語を含む。コンテンツ情報格納部 3 2 0 は、検索履歴を格納するとともに、コンテンツが検索された頻度を検索語毎に集計し、その集計結果を格納する。分類部 3 4 0 は、分類情報格納部 3 3 0 に格納された情報と、コンテンツ情報格納部 3 2 0 に格納された情報とを用いて、カテゴリが未知のコンテンツを分類する。分類部 3 4 0 は、カテゴリが未知のコンテンツの検索に用いられた検索語とカテゴリが既知のコンテンツの検索に用いられた検索語（及びそのカテゴリ）とに基づいて、カテゴリが未知のコンテンツを分類する。

【選択図】 図 3



【特許請求の範囲】**【請求項 1】**

所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得する第 1 の取得部と、

前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得する第 2 の取得部と、

前記第 1 の取得部により取得された検索履歴と前記第 2 の取得部により取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類する分類部と

を備える分類装置。

10

【請求項 2】

前記分類部は、前記検索語毎の検索された頻度に基づいて前記第 2 のコンテンツを分類する

ことを特徴とする請求項 1 に記載の分類装置。

【請求項 3】

前記分類部は、検索された頻度に応じたスコアを前記検索語毎に算出し、当該スコアを用いて前記第 2 のコンテンツを分類する

ことを特徴とする請求項 2 に記載の分類装置。

【請求項 4】

URL (Uniform Resource Locator) の少なくとも一部が共通する前記第 1 のコンテンツ又は前記第 2 のコンテンツどうしが同一のカテゴリに分類されることを特徴とする請求項 1 ないし 3 のいずれかに記載の分類装置。

20

【請求項 5】

請求項 1 ないし 4 のいずれかに記載の分類装置と、ユーザにより入力された検索語に対応するコンテンツを当該ユーザに提示する検索装置とを有し、

前記検索装置が、前記入力された検索語に対応するコンテンツを前記カテゴリ毎に分類した態様で提示するための提示部を備える

ことを特徴とするコンテンツ検索システム。

【請求項 6】

所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、

30

前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得するステップと、

前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類するステップと

を有することを特徴とするコンテンツ分類方法。

【請求項 7】

所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、

前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得するステップと、

40

前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類するステップと、

入力された検索語に対応する前記第 1 のコンテンツ又は前記第 2 のコンテンツを前記カテゴリ毎に分類した態様でユーザに提示するステップと

を有することを特徴とするコンテンツ検索方法。

【請求項 8】

コンピュータに、

所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、

50

前記カテゴリが未知である第2のコンテンツを検索するために用いられた検索語を取得するステップと、

前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第2のコンテンツを前記カテゴリに従って分類するステップと

を実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンテンツを分類し、又は検索するための技術に関する。

【背景技術】

【0002】

インターネットを介して所望のコンテンツを得るために、サーチエンジン（検索エンジン）が利用されている。Webページの検索は、例えば、ユーザが検索語（キーワード）を入力し、その検索語を文字列として含んでいるWebページをユーザに検索結果として提示する、といった手順で行われる。なお、検索語に対応するWebページが複数ある場合には、サーチエンジン毎の規則やアルゴリズムに従って順位付けが行われる。このような検索手法は、「キーワード検索」と呼ばれている。また、キーワード検索のほかにも、あらかじめ設定されたカテゴリに従ってWebサイトやWebページを検索する、いわゆる「カテゴリ検索」も知られている。

【0003】

特許文献1には、検索に必要なキーワードを抽出するための技術が記載されている。また、特許文献2、3には、形態素解析等の自然言語処理を用いて文書を分類するための技術が記載されている。

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2002-149683号公報

【特許文献2】特開平2-158871号公報

【特許文献3】特開平11-328211号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

ところで、コンテンツを検索する場合には、1語や2語の検索語で検索することも多い。この傾向は、ユーザが文字入力に不慣れな場合や、文字入力に特化していないデバイス（例えば、携帯電話機など）で文字を入力する場合に、より顕著である。検索語そのものの情報量が少ない場合には、ユーザが意図していないコンテンツまでもが検索結果に含まれてしまうことも多い。

【0006】

また、カテゴリ検索には、あらかじめ分類されているコンテンツでなければ検索結果として利用できないという問題があり、コンテンツの頻繁な追加や更新に対応することが困難である。さらに、自然言語処理のように、コンテンツの内容そのものを解析する場合には、その処理に時間を要するだけでなく、一つ一つのコンテンツの情報量（文字数）が増えるほど処理時間も増加してしまう。

【0007】

そこで、本発明は、コンテンツを、当該コンテンツの内容によらないで分類できるようにすることを目的とする。

【課題を解決するための手段】

【0008】

本発明の一態様に係る分類装置は、所定のカテゴリに従って分類された第1のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得

10

20

30

40

50

する第 1 の取得部と、前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得する第 2 の取得部と、前記第 1 の取得部により取得された検索履歴と前記第 2 の取得部により取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類する分類部とを備える。

【 0 0 0 9 】

好ましい態様において、前記分類部は、前記検索語毎の検索された頻度に基づいて前記第 2 のコンテンツ进行分类する。

さらに好ましい態様において、前記分類部は、検索された頻度に応じたスコアを前記検索語毎に算出し、当該スコアを用いて前記第 2 のコンテンツ进行分类する。

他の好ましい態様において、前記分類装置は、U R L (Uniform Resource Locator) の少なくとも一部が共通する前記第 1 のコンテンツ又は前記第 2 のコンテンツどうしが同一のカテゴリに分類されることを特徴とする。

【 0 0 1 0 】

本発明の他の態様に係るコンテンツ検索システムは、前記分類装置と、ユーザにより入力された検索語に対応するコンテンツを当該ユーザに提示する検索装置とを有し、前記検索装置が、前記入力された検索語に対応するコンテンツを前記カテゴリ毎に分類した態様で提示するための提示部を備える。

【 0 0 1 1 】

本発明の他の態様に係るコンテンツ分類方法は、所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得するステップと、前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類するステップとを有する。

【 0 0 1 2 】

本発明の他の態様に係るコンテンツ検索方法は、所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得するステップと、前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類するステップと、入力された検索語に対応する前記第 1 のコンテンツ又は前記第 2 のコンテンツを前記カテゴリ毎に分類した態様でユーザに提示するステップとを有する。

【 0 0 1 3 】

本発明の他の態様に係るプログラムは、コンピュータに、所定のカテゴリに従って分類された第 1 のコンテンツと、当該コンテンツを検索するために用いられた検索語とを対応付けた検索履歴を取得するステップと、前記カテゴリが未知である第 2 のコンテンツを検索するために用いられた検索語を取得するステップと、前記取得された検索履歴と前記取得された検索語とを比較することによって、前記第 2 のコンテンツを前記カテゴリに従って分類するステップとを実行させるためのものである。

【 発明の効果 】

【 0 0 1 4 】

本発明によれば、コンテンツを、当該コンテンツの内容によらないで分類することが可能である。

【 図面の簡単な説明 】

【 0 0 1 5 】

【 図 1 】 本発明の一実施形態の全体構成を示す図

【 図 2 】 コンテンツ検索システムのハードウェア構成を示すブロック図

【 図 3 】 分類装置の機能的構成を示す機能ブロック図

【 図 4 】 検索履歴情報のデータ構造を例示する図

【 図 5 】 頻度情報のデータ構造を例示する図

【図 6】分類情報のデータ構造を例示する図

【図 7】分類装置によるコンテンツの分類方法を示すフローチャート

【図 8】検索装置によるコンテンツの検索方法を示すフローチャート

【図 9】分類装置及び記憶装置の構成（変形例）を示すブロック図

【発明を実施するための形態】

【0016】

[実施形態]

図 1 は、本発明の一実施形態の全体構成を示す図である。本実施形態のコンテンツ検索システム 110 は、ユーザがコンテンツを検索するために用いられるものであり、複数のユーザのクライアント端末 120 によってアクセスされる。クライアント端末 120 は、通信ネットワーク 130 を介してコンテンツ検索システム 110 にアクセスするコンピュータ装置であり、例えば、パーソナルコンピュータ、携帯電話機、スマートフォンなどである。クライアント端末 120 は、通信ネットワーク 130 と通信を行う手段と、文字入力等の操作を受け付ける手段と、検索結果を表示する手段とを少なくとも備える。通信ネットワーク 130 は、インターネット、イントラネット、移動体通信網などであり、また、これらを組み合わせた複合的なネットワークであってもよい。

【0017】

本実施形態において、コンテンツとは、ユーザが視覚的に閲覧可能な情報を含むひとまとまりのデータをいい、ここでは、文字コード（ASCII、Unicode、Shift_JIS等）によって記述された文字列を少なくとも含むものとする。コンテンツは、例えば、HTML（HyperText Markup Language）形式のWebページであるが、PDF（Portable Document Format）データをはじめ、通信ネットワーク 130 を介してやりとりが可能なさまざまなデータを含み得る。ここでいうWebページは、あらかじめ記憶されているものに限らず、CGI（Common Gateway Interface）などによって動的に生成されたものであってもよい。コンテンツは、通信ネットワーク 130 に接続された図示せぬWebサーバに記憶されている。また、コンテンツは、音楽データのような視覚的なデータ以外のデータを含んでいてもよい。

【0018】

コンテンツ検索システム 110 は、検索装置 111 と、分類装置 112 とを備える。検索装置 111 は、サーチエンジンの機能を有し、クライアント端末 120 から受け付けた検索クエリに応じた検索結果をクライアント端末 120 に提示する。検索クエリは、コンテンツを検索するために用いられる 1 又は複数の検索語を少なくとも含む。また、検索クエリは、クライアント端末 120（又はそのユーザ）を識別する情報や、検索語以外の検索条件（検索の態様を指定する条件。前方一致検索、アンド検索など。）を含んでもよい。検索結果は、検索語に対応するコンテンツにリンクするURLを含んだWebページによってクライアント端末 120 に提供される。

【0019】

分類装置 112 は、クライアント端末 120 のユーザによって検索されたコンテンツを当該コンテンツの特徴によって分類する機能を有する。分類装置 112 による分類は、あらかじめ決められたカテゴリを用いて行われる。本実施形態におけるカテゴリは、例えば、「スポーツ」、「映画」、「ギャンブル」、「音楽」、「占い」、「グルメ」といったものである。なお、カテゴリは、コンテンツ検索システム 110 の利用目的や利用するユーザに応じて適宜定められればよい。

【0020】

図 2 は、コンテンツ検索システム 110 のハードウェア構成を示すブロック図である。検索装置 111 は、図 2 に示すように、制御部 211 と、記憶部 212 と、第 1 通信部 213 と、第 2 通信部 214 とを備える。制御部 211 は、検索装置 111 の各部の動作を制御する手段である。制御部 211 は、CPU（Central Processing Unit）等の演算処理装置や主記憶装置に相当する記憶手段（メインメモリ）を備え、プログラムを実行することによって検索機能を実現する。ここにおいて、検索機能とは、検索語に応じた検索結

果を提示する機能をいい、本発明に係る提示部に相当するものである。記憶部 2 1 2 は、H D D (Hard Disk Drive) 等の補助記憶装置に相当する記憶手段を備え、検索機能に必要なデータ (例えば、いわゆるインデクス等) を記憶する。第 1 通信部 2 1 3 は、分類装置 1 1 2 と通信するためのインターフェースを備える。第 2 通信部 2 1 4 は、通信ネットワーク 1 3 0 と通信するためのインターフェースである。

【0021】

分類装置 1 1 2 は、図 2 に示すように、制御部 2 2 1 と、記憶部 2 2 2 と、第 1 通信部 2 2 3 とを備える。これらの各部は、検索装置 1 1 1 の同名の構成要素と同様のハードウェア構成を有する。ただし、記憶部 2 2 2 は、その記憶するデータが記憶部 2 1 2 とは相違し、コンテンツの分類に必要なデータを記憶している。

10

【0022】

図 3 は、分類装置 1 1 2 の機能的構成を示す機能ブロック図である。分類装置 1 1 2 は、プログラムを実行することにより、図 3 に示す履歴受信部 3 1 0、コンテンツ情報格納部 3 2 0、分類情報格納部 3 3 0 及び分類部 3 4 0 の各部に相当する機能を実現する。また、コンテンツ情報格納部 3 2 0 は、より詳細には、検索履歴格納部 3 2 1、検索履歴集計部 3 2 2 及び頻度格納部 3 2 3 に機能的に分類される。

【0023】

履歴受信部 3 1 0 は、検索装置 1 1 1 から検索履歴を受信する。ここにおいて、検索履歴とは、検索語と、その検索語を用いて検索され、ユーザが閲覧するためにクライアント端末 1 2 0 で選択したコンテンツとを対応付けたものをいう。ここでいう「選択」とは、ユーザの操作 (クリック等) に基づくものである。履歴受信部 3 1 0 により受信されるコンテンツは、分類の対象であるコンテンツ、すなわち、まだ分類されておらず、カテゴリが未知であるコンテンツと、カテゴリが既知であるコンテンツの双方が含まれ得る。

20

【0024】

コンテンツ情報格納部 3 2 0 は、記憶部 2 2 2 にコンテンツ情報を格納する。コンテンツ情報は、検索履歴情報と頻度情報の総称である。検索履歴情報は、検索履歴を表すデータである。検索履歴格納部 3 2 1 は、履歴受信部 3 1 0 により受信された検索履歴に基づき、記憶部 2 2 2 に検索履歴情報を格納する。

【0025】

頻度情報は、あるコンテンツがどの検索語によって何回検索されたかを表すデータである。本実施形態において、頻度情報は、コンテンツを表す URL と、当該 URL に対応する検索語と、その検索された頻度とを対応付けたデータである。なお、ここでいう「頻度」は、検索クエリとして得られた延べ回数ではなく、例えば、同一のユーザが同一の検索語で同一のコンテンツを繰り返し何回も検索した場合には、これらを 1 回の検索とみなした値であってもよい。この場合、頻度の値は、単位時間当たりにあるコンテンツをある検索語で検索したユーザの人数に相当する。

30

【0026】

頻度情報は、検索履歴集計部 3 2 2 によって集計され、頻度格納部 3 2 3 によって記憶部 2 2 2 に格納される。検索履歴集計部 3 2 2 は、検索履歴格納部 3 2 1 によって格納された検索履歴情報を読み出し、その検索履歴情報の中で重複しているコンテンツを検索語毎に集計し、頻度を算出する。頻度格納部 3 2 3 は、検索履歴集計部 3 2 2 による集計結果を記憶部 2 2 2 に格納する。

40

【0027】

図 4 は、検索履歴情報のデータ構造を例示する図である。この例において、検索履歴情報は、ユーザ ID と、タイムスタンプと、検索語と、URL とを対応付けて記述したデータである。ユーザ ID は、ユーザを一意的に識別するためのデータである。ユーザの識別は、パスワード等による認証や Cookie の利用など、周知の適当な技術によって実現されればよい。タイムスタンプは、ユーザが検索を行った時刻 (以下「検索時刻」という。) を表すデータである。検索時刻は、クライアント端末 1 2 0 で検索クエリの送信時に計測されてもよいし、検索クエリの受信時に検索装置 1 1 1 で計測されてもよい。

50

図 5 は、頻度情報のデータ構造を例示する図である。この例において、頻度情報は、URL と、検索語と、頻度とを対応付けて記述したデータである。

なお、検索履歴情報又は頻度情報は、コンテンツが分類済みであるか否か（すなわち、カテゴリが既知か未知か）を表すフラグを含んでいてもよい。

【0028】

分類情報格納部 330 は、コンテンツの分類結果を表す分類情報を記憶部 222 に格納する。分類情報は、記憶部 222 にあらかじめ記憶されているものと、分類部 340 による分類結果として事後的に格納されるものとがある。本実施形態の分類情報は、URL と、カテゴリ ID と、スコアとを対応付けて記述したデータである。カテゴリ ID は、既知のカテゴリを一意的に識別するために割り当てられるデータである。また、スコアは、URL によって表されるコンテンツがそれぞれのカテゴリに適合する度合いを示す値であり、本実施形態においては、「1」が最大（最も適合する）で「0」が最小であるとする。

【0029】

図 6 は、分類情報のデータ構造を例示する図である。この例においては、URL 「http://aaa.com/aaa.cgi?a1=1&a2=2」で表されるコンテンツは、カテゴリ ID 「24」で表されるカテゴリに最も適合している（相応しい）ということになる。なお、カテゴリは、一つのコンテンツにつき 1 種類でなくてもよい。例えば、コンテンツは、所定のスコア以上のカテゴリのすべてに属するとしてもよいし、スコアの大きい順に所定数のカテゴリに属するとしてもよい。また、図 6 においては、スコアが「0」である分類情報の表示を省略しているが、実際のデータは、各 URL と各カテゴリによって考えられる組み合わせのそれぞれについて、「0」かそれ以外のスコアが付与されているものとする。

【0030】

分類部 340 は、分類情報格納部 330 に格納された分類情報を用いて、カテゴリが未知であるコンテンツを分類し、その分類結果を分類情報格納部 330 に供給する。分類部 340 は、カテゴリが既知のコンテンツの検索履歴とカテゴリが未知のコンテンツを検索するために用いられた検索語とを比較することによって、カテゴリが未知のコンテンツを分類する。

【0031】

本実施形態において、分類部 340 は、コンテンツ情報格納部 320 に格納されたコンテンツ情報からカテゴリが未知であるコンテンツに関するデータを取得し、これを分類の対象とする。次に、分類部 340 は、分類情報格納部 330 を介して分類情報を取得し、分類モデルを生成する。ここにおいて、分類モデルとは、コンテンツがどのカテゴリに属するのが適当であるかを数学的にモデル化したものである。分類部 340 は、分類モデルを用いて、分類対象のコンテンツのそれぞれにどのカテゴリが適しているかを算出し、その算出結果を分類情報として分類情報格納部 330 に供給する。分類モデルは、例えば、SVM (Support Vector Machine) やニューラルネットワークを用いて実現可能である。

【0032】

以上の構成のもと、コンテンツ検索システム 110 は、クライアント端末 120 からの要求（すなわち検索クエリ）に応じて検索を実行し、検索結果をクライアント端末 120 に提示する。また、コンテンツ検索システム 110 は、適当なタイミングで、カテゴリが未知のコンテンツに関する情報を取得し、コンテンツの分類を繰り返し実行する。

また、コンテンツ検索システム 110 は、周知のキーワード検索による検索結果と、コンテンツをカテゴリ毎に分類して表示する検索結果とをユーザに提供することができる。コンテンツ検索システム 110 は、これらの検索結果を両方提示してもよいし、ユーザによる事前の設定に応じていずれかの検索結果を択一的に提示するようにしてもよい。

【0033】

図 7 は、分類装置 112 によるコンテンツの分類方法を示すフローチャートである。図 7 に示すように、分類装置 112 の制御部 221 は、最初に既知の分類情報を格納する（ステップ S11）。ステップ S11 において、制御部 221 は、カテゴリに関する既存の情報源を利用して分類情報を得る。かかる情報源としては、例えば、ODP (Open Direc

10

20

30

40

50

tory Project)に登録されているWebサイトを利用可能である。あるいは、周知の代表的なサーチエンジンを用いてカテゴリ名によって検索を実行し、その検索結果の上位に提示されるWebサイトを分類情報に利用したり、既存のポータルサイトによって分類されているカテゴリを分類情報に利用したりすることも可能である。なお、分類装置112が用いるカテゴリと既存の情報源によるカテゴリとに相違がある場合には、あらかじめ対応表を作成しておき、そのときに用いるカテゴリに置き換えられるようにすることが望ましい。また、分類情報のうち、スコアは、適当な値があらかじめ機械的に(あるいは手作業で)設定される。

【0034】

次に、制御部221は、検索履歴を検索装置111から受信し、これを検索履歴情報として格納する(ステップS12)。制御部221は、検索クエリが発生する毎に検索履歴を受信してもよいし、検索履歴を所定時間毎又は所定数毎に一括して受信する動作を繰り返してもよい。このような動作が繰り返されることにより、検索履歴情報が徐々に蓄積される。一定量の検索履歴情報が蓄積されると、制御部221は、頻度情報をコンテンツ毎に集計し、格納する(ステップS13)。

【0035】

続いて、制御部221は、このようにして格納された分類情報と頻度情報とに基づき、分類モデルを生成する(ステップS14)。分類モデルの具体的な生成方法の一例は、SVMによる回帰分析を行う場合であれば、以下のとおりである。

【0036】

制御部221は、カテゴリが既知であるコンテンツの特徴量として、頻度情報をベクトル形式で表現した X_i を用いる。 X_i は、具体的には、以下の(1)式で表される。ここにおいて、 x_{ij} は、各コンテンツ及び検索語を一意的に識別するためのIDをそれぞれ i 、 j とした場合の頻度の値であり、コンテンツの総数を N 、検索語の総数を V とすると、 $1 \leq i \leq N$ 、 $1 \leq j \leq V$ をそれぞれ満たす。つまり、 X_i は、コンテンツの検索語毎の頻度をコンテンツ毎に表すものである。

$$X_i = (x_{i1}, x_{i2}, x_{i3}, \dots, x_{iV}) \quad \dots (1)$$

【0037】

なお、制御部221は、必要に応じて、 X_i に対数化や正規化を行ってもよい。対数化や正規化は、 x_{ij} のそれぞれの差が大きすぎ、大小関係の特徴が強く現れすぎる場合に有効である。対数化は、例えば、 x_{ij} を $\log_{10}(x_{ij} + 1)$ 、すなわち10を底とする($x_{ij} + 1$)の対数に置き換えることで実現可能である。また、正規化は、 X_i の絶対値が1となるように、各要素(x_{ij})を要素の二乗和で除算することで実現可能である。

【0038】

また、制御部221は、カテゴリの特徴量として、分類情報をベクトル形式で表現した Y_i を用いる。 Y_i は、具体的には、以下の(2)式で表される。ここにおいて、 y_{ik} は、コンテンツのIDを i 、カテゴリIDを k とした場合のスコアである。カテゴリIDは、カテゴリの総数を L とすると、 $1 \leq k \leq L$ を満たす。つまり、 Y_i は、コンテンツのカテゴリ毎のスコアをコンテンツ毎に表すものである。

$$Y_i = (y_{i1}, y_{i2}, y_{i3}, \dots, y_{iL}) \quad \dots (2)$$

【0039】

制御部221は、これらの特徴量から分類モデルを生成する。分類モデルは、 X_i を説明変数、 Y_i を目的変数とするものであり、以下においては $F_k(X_i)$ と表現する。 $F_k(X_i)$ は、カテゴリ毎(すなわちカテゴリID毎)に生成される。制御部221は、このようにして生成される分類モデル($F_k(X_i)$)に対して、カテゴリが未知であるコンテンツの特徴量である X_i' を入力する(なお、 X_i と X_i' の相違点は、カテゴリが既知であるか未知であるかの1点のみである)。そうすると、制御部221は、カテゴリが未知であるコンテンツの目的変数として、上述した Y_i に相当する特徴量を得ることができる。

【0040】

10

20

30

40

50

制御部 2 2 1 は、このようにして生成された分類モデルを用いて、カテゴリが未知であるコンテンツについてカテゴリ毎のスコアを算出し、その算出結果を分類情報として格納する（ステップ S 1 5）。なお、制御部 2 2 1 は、スコアに対して閾値を設定し、閾値を下回るスコアをすべて「0」とであるとみなしてもよい。このようにすることで、実態に即していない分類がされることを防ぐことが可能である。また、制御部 2 2 1 は、タイムスタンプに基づいて頻度又はスコアに対して重み付けを行い、より新しい検索履歴ほど分類に強い影響を与えるようにしてもよい。

【0041】

制御部 2 2 1 は、一定量の検索履歴が新たに蓄積される毎に、あるいは一定間隔で、図 7 に示す分類処理を繰り返す。ただし、制御部 2 2 1 は、いったん分類情報が格納された後には、ステップ S 1 1 の処理をスキップし、ステップ S 1 2 の処理から実行すればよい。また、ステップ S 1 5 の処理によって分類情報が新たに格納されたコンテンツは、その後はカテゴリが既知のコンテンツとして扱われる。

【0042】

図 8 は、検索装置 1 1 1 によるコンテンツの検索方法を示すフローチャートである。図 8 に示すように、検索装置 1 1 1 の制御部 2 1 1 は、ユーザから検索クエリを受信することにより、検索語を取得する（ステップ S 2 1）。制御部 2 1 1 は、取得した検索語に基づき、周知のキーワード検索によってコンテンツのリストを生成する（ステップ S 2 2）。ステップ S 2 2 において生成されるリストは、検索語を含むコンテンツを適当に順位付けしたものである。

【0043】

次に、制御部 2 1 1 は、分類装置 1 1 2 に格納されている分類情報を参照することにより、ステップ S 2 2 において生成されたリストに含まれるコンテンツをカテゴリ毎に分類する（ステップ S 2 3）。そして、制御部 2 1 1 は、ユーザが入力した検索語に対応するコンテンツをカテゴリ毎に分類した態様で表示するための検索結果情報を当該ユーザのクライアント端末 1 2 0 に送信する（ステップ S 2 4）。検索結果情報の表示態様は、例えば、各カテゴリのコンテンツのリストをスコアが高い順に表示するものであってもよいし、各カテゴリのコンテンツのリストをタブで切り替えて表示するものであってもよい。なお、制御部 2 1 1 は、コンテンツに対応する要約文を検索結果情報に含めてもよい。このようにすれば、ユーザによる各コンテンツの取捨選択を容易にすることができる。また、制御部 2 1 1 は、カテゴリが未知のコンテンツが検索結果に含まれる場合などには、必要に応じて、キーワード検索を行っただけの（未分類の）リストを検索結果情報に含めてもよい。

【0044】

検索結果情報には、目的のコンテンツに直接リンクする URL ではなく、いったんコンテンツ検索システム 1 1 0 を経由して目的のコンテンツにリンクする URL（いわゆるリダイレクト URL）が記述されている。この URL には、パラメータとして、目的のコンテンツの URL に加え、ユーザが入力した検索語が含まれる。このようにすることで、コンテンツ検索システム 1 1 0 は、コンテンツを検索するために用いられた検索語を特定し、検索履歴を取得することが可能である。

【0045】

以上のとおり、本実施形態によれば、コンテンツを検索するために用いられた検索語を利用することで、コンテンツの内容によらない分類を行うことが可能となる。よって、本実施形態によれば、コンテンツに対して自然言語処理を実行することなくコンテンツを分類することが可能となる。ただし、コンテンツの内容によらない分類が有効に機能するためには、当該コンテンツの検索が既にある程度行われていることが条件となる。

【0046】

また、本実施形態によれば、コンテンツを検索するために実際に用いられた検索語を利用することで、各ユーザの意図に即した分類を行うことが可能であるともいえる。例えば、携帯電話機向けのコンテンツにおいては、画面サイズ等の表示上の制約から、文字の情

10

20

30

40

50

報量が比較的少ない場合がある。かかるコンテンツをその内容（すなわち、コンテンツに含まれる文字列）に基づいて分類した場合、その分類結果にユーザの意図を的確に反映させられず、分類の精度・確度が低下する可能性がある。よって、コンテンツ検索システム 110 は、文字列としての情報量が比較的少ないコンテンツの検索に適用するのに好適であるといえる。

【0047】

[変形例]

上述した実施形態は、本発明の実施の一態様である。本発明は、上述した実施形態に対して以下の変形を適用した態様で実施することも可能である。なお、以下に示す変形例は、必要に応じて、各々を適当に組み合わせて実施されてもよいものである。

【0048】

(変形例 1)

上述したとおり、本発明に係るコンテンツは、その分類に際して自然言語処理が不要である。すなわち、本発明に係るコンテンツは、文字列を含まないデータであっても分類可能である。したがって、本発明は、文字列を含むか否かを問わず、ユーザが検索可能なあらゆるデータを分類の対象にすることができる。

【0049】

(変形例 2)

本発明に係る分類装置は、URL の少なくとも一部が共通する複数のコンテンツがある場合に、これらが同一のカテゴリに分類されるように動作するものであってもよい。このようにすれば、カテゴリが未知のコンテンツをより高速に分類することが可能であるとともに、検索されたことが一度もないコンテンツであっても分類が可能になる場合がある。なお、本例において、カテゴリの共通化は、ドメイン名、ホスト名、FQDN (Fully Qualified Domain Name) などを単位として行われる。また、URL がパス名 (ディレクトリ名) によって階層分けされている場合には、かかる階層がカテゴリに対応付けられていてもよい。さらに、URL にパラメータ部 (図 4 等の例における「?」以降の文字列) が含まれる場合には、パラメータ部以外の部分が共通している URL のカテゴリを共通化することも可能である。

【0050】

(変形例 3)

上述した実施形態においては、コンテンツの分類処理は、カテゴリが未知のコンテンツに対してのみ実行された。しかし、本発明は、カテゴリが既知のコンテンツにも分類処理を実行し、必要に応じて、コンテンツのカテゴリを変更できるようにしてもよい。このようにすれば、コンテンツの内容が時間の経過に応じて (URL は変えずに) 変更されたとしても、より適切なカテゴリに再分類することが可能となる。具体的には、制御部 221 は、カテゴリが既知のコンテンツに対して上述した分類モデルを適用し、既に算出されたスコアと新たに算出されたスコアとの間に一定以上の乖離がある場合に、既に算出されたスコアを新たに算出されたスコアに書き換えるようにしてもよい。

【0051】

(変形例 4)

本発明は、コンテンツ情報や分類情報の記憶手段を分類装置から分離した態様でも実施可能である。すなわち、本発明に係る分類装置は、コンテンツ情報や分類情報を格納する手段を別体に構成したものであってもよい。

【0052】

図 9 は、本例に係る分類装置及び記憶装置の構成を示すブロック図である。図 9 に示すように、分類装置 910 は、第 1 取得部 911 と、第 2 取得部 912 と、分類部 913 と、供給部 914 とを備える。また、記憶装置 920 は、コンテンツ情報を記憶する手段であり、上述した実施形態のコンテンツ情報格納部 320 に相当する機能を少なくとも備える。記憶装置 930 は、分類情報を記憶する手段であり、上述した実施形態の分類情報格納部 330 に相当する機能を少なくとも備える。なお、記憶装置 920 は、分類装置 91

10

20

30

40

50

0を実施する事業者とは異なる事業者のサーチエンジンであってもよい。

【0053】

第1取得部911は、記憶装置920からコンテンツ情報（検索履歴情報、頻度情報）を取得する。第2取得部912は、記憶装置930から分類情報を取得する。分類部913は、上述した実施形態の分類部340に相当する機能を備える。供給部914は、分類部913による分類結果を記憶装置930に供給する。

【0054】

（変形例5）

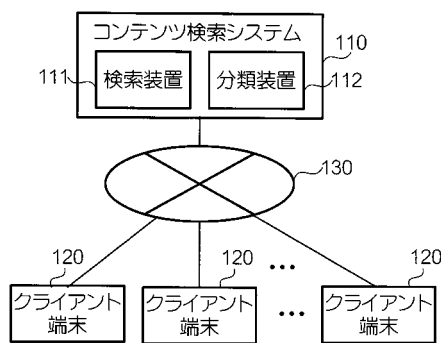
本発明は、上述した実施形態のように分類装置と検索装置とを別体にするのではなく、これらを単一のコンピュータ装置で実現することも可能である。また、本発明は、コンテンツを分類し、又は検索するための方法や、コンピュータ装置を上述した分類装置や検索装置として機能させるためのプログラムとしても提供可能である。かかるプログラムは、光ディスク等の記録媒体に記録した形態で提供されたり、インターネット等のネットワークを介して、コンピュータにダウンロードさせ、これをインストールして利用可能にするなどの形態で提供されたりすることも可能である。

【符号の説明】

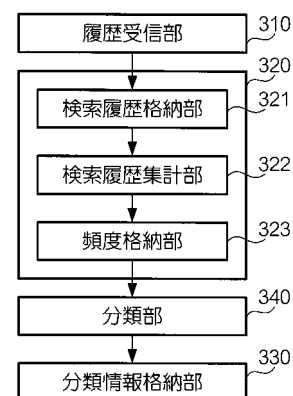
【0055】

110...コンテンツ検索システム、111...検索装置、112...分類装置、120...クライアント端末、130...通信ネットワーク、211、221...制御部、212、222...記憶部、213、223...第1通信部、214...第2通信部、310...履歴受信部、320...コンテンツ情報格納部、321...検索履歴格納部、322...検索履歴集計部、323...頻度格納部、330...分類情報格納部、340...分類部、910...分類装置、911...第1取得部、912...第2取得部、913...分類部、914...供給部、920、930...記憶装置

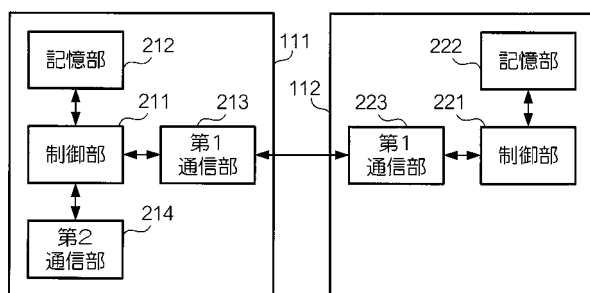
【図1】



【図3】



【図2】



【図 4】

ユーザ ID	タイムスタンプ	検索語	URL
000	2010/5/28 07:00:00	A	http://aaa.com/aaa.cgi?a1=1&a2=2
000	2010/5/28 07:31:15	B	http://bbb.com/bbb.html
000	2010/5/28 08:02:30	C	http://bbb.com/bbb.cgi?b1=1
001	2010/5/28 09:35:10	D	http://ddd.co.jp/ddd/ddd.gif
001	2010/5/28 12:40:02	E	http://eee.com/eee.cgi?e1=1;sessionId=111
002	2010/5/28 08:24:10	F	http://fff.com/fff.html
002	2010/5/28 08:27:28	G	http://fff.com/fff/fff/fff.cgi?id=fff
⋮	⋮	⋮	⋮

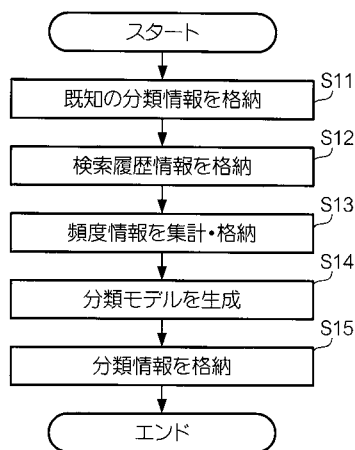
【図 6】

URL	カテゴリID	スコア
http://aaa.com/aaa.cgi?a1=1&a2=2	1	0.45
http://aaa.com/aaa.cgi?a1=1&a2=2	5	0.81
http://aaa.com/aaa.cgi?a1=1&a2=2	24	1.00
http://bbb.com/bbb.html	3	0.26
http://bbb.com/bbb.html	9	0.71
http://bbb.com/bbb.html	11	0.57
http://bbb.com/bbb.html	41	0.98
⋮	⋮	⋮

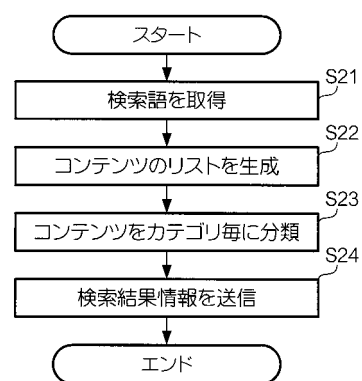
【図 5】

URL	検索語	頻度
http://aaa.com/aaa.cgi?a1=1&a2=2	A	34
http://aaa.com/aaa.cgi?a1=1&a2=2	C	664
http://aaa.com/aaa.cgi?a1=1&a2=2	D	65
http://aaa.com/aaa.cgi?a1=1&a2=2	F	216
http://aaa.com/aaa.cgi?a1=1&a2=2	H	45
http://bbb.com/bbb.html	B	74
http://bbb.com/bbb.html	C	57
⋮	⋮	⋮

【図 7】



【図 8】



【 図 9 】

