

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7434867号
(P7434867)

(45)発行日 令和6年2月21日(2024.2.21)

(24)登録日 令和6年2月13日(2024.2.13)

(51)国際特許分類	F I
G 0 6 F 16/957 (2019.01)	G 0 6 F 16/957
G 0 6 F 16/951 (2019.01)	G 0 6 F 16/951
H 0 4 L 67/02 (2022.01)	H 0 4 L 67/02

請求項の数 10 (全16頁)

(21)出願番号	特願2019-223095(P2019-223095)	(73)特許権者	000005223 富士通株式会社
(22)出願日	令和1年12月10日(2019.12.10)		神奈川県川崎市中原区上小田中4丁目1番1号
(65)公開番号	特開2020-98596(P2020-98596A)	(74)代理人	100107766 弁理士 伊東 忠重
(43)公開日	令和2年6月25日(2020.6.25)	(74)代理人	100070150 弁理士 伊東 忠彦
審査請求日	令和4年8月9日(2022.8.9)	(72)発明者	ジョン・ジョングアン 中国, 1 0 0 0 2 7, ベイジン, チャオヤン ディストリクト, ゴン ティ ベイ ルウ ナンバー 2 エイ, パシフィック センチュリー プレイス, スペース 8, ゲート 6, ユニット 3 エフ, 3 5 5 富士通研究開発中心有限公司内
(31)優先権主張番号	201811549030.2		最終頁に続く
(32)優先日	平成30年12月18日(2018.12.18)		
(33)優先権主張国・地域又は機関	中国(CN)		

(54)【発明の名称】 ウェブページから情報を抽出する方法、装置及び記憶媒体

(57)【特許請求の範囲】

【請求項1】

ウェブページから情報を抽出する方法であって、コンピュータが、
前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成するステップと、
前記木構造におけるナビゲーションバーノードを決定するステップと、
前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するステップと、
マッチする葉ノードに対応するページにおける情報を抽出するステップと、を実行する方法。

【請求項2】

前記木構造におけるナビゲーションバーノードを決定するステップは、
前記木構造に出現する回数が所定閾値よりも大きい葉ノードのみを含む非葉ノードを決定するステップと、
前記非葉ノードを並び替えて前記ナビゲーションバーノードを決定するステップと、を含む、請求項1に記載の方法。

【請求項3】

葉ノードの出現回数が所定閾値よりも大きいか否かを決定することは、
前記葉ノードのテキスト及び経路情報の前記木構造における出現回数が前記所定閾値よりも大きいか否かを決定すること、を含む、請求項2に記載の方法。

10

20

【請求項 4】

前記経路情報は、前記葉ノードからその n 番目の先祖ノードまでの経路であり、n は正整数である、請求項 3 に記載の方法。

【請求項 5】

n は 5 以上である、請求項 4 に記載の方法。

【請求項 6】

前記非葉ノードを並び替えて前記ナビゲーションバーノードを決定するステップは、前記非葉ノードの特徴値を計算するステップであって、前記特徴値は、前記非葉ノードによりカバーされる葉ノードの数及び前記回数により決定される、ステップと、前記非葉ノードのうちの最大の特徴値を有する非葉ノードを前記ナビゲーションバーノードとして決定するステップと、を含む、請求項 2 に記載の方法。

10

【請求項 7】

マッチする葉ノードに対応するページにおける情報を抽出するステップは、前記マッチする葉ノードに対応するページに含まれるターゲットノードを決定するステップと、前記ターゲットノードによりカバーされる各葉ノードのテキストをそれぞれ抽出するステップと、を含む、請求項 1 乃至 6 の何れかに記載の方法。

【請求項 8】

前記ターゲットノードは、前記ターゲットノードに含まれる各葉ノードのテキスト及び経路情報の前記木構造における出現回数が所定閾値以下であること、前記ターゲットノードが、前記木構造に出現する回数が所定閾値よりも大きい葉ノードのみを含む非葉ノードのうちの非葉ノードではないこと、及び前記ターゲットノードに含まれる全ての葉ノードのテキストの合計長さが該木構造における他の非葉ノードのテキストの合計長さよりも大きいこと、により決定される、請求項 7 に記載の方法。

20

【請求項 9】

ウェブページから情報を抽出する装置であって、前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成する木構造生成手段と、前記木構造におけるナビゲーションバーノードを決定するナビゲーションバーノード決定手段と、前記ナビゲーションバーノードによりカバーされる、1 つ又は複数のキーワードにマッチする葉ノードを決定するマッチノード決定手段と、マッチする葉ノードに対応するページにおける情報を抽出する情報抽出手段と、を含む、装置。

30

【請求項 10】

プログラムを記憶したコンピュータ読み取り可能な記憶媒体であって、前記プログラムがプロセッサにより実行される際に、ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成するステップと、前記木構造におけるナビゲーションバーノードを決定するステップと、前記ナビゲーションバーノードによりカバーされる、1 つ又は複数のキーワードにマッチする葉ノードを決定するステップと、マッチする葉ノードに対応するページにおける情報を抽出するステップと、を実行させる、記憶媒体。

40

【発明の詳細な説明】**【技術分野】****【0001】**

本開示は、自然言語処理に関し、具体的には、複数のウェブページに基づく情報抽出に

50

関する。

【背景技術】

【0002】

インターネットから情報を収集、抽出することは、知識ベースを構築する重要な手段である。例えば、電子商取引会社のウェブページから製品の情報を抽出し、製品の知識ベースを構築することができる。従来の方法は、主に2種類の方法に分類される。

【0003】

1種類目の方法は、類似の構造を有するページ（例えば、電子商取引会社のウェブサイトの製品リストのページでは、各ページの構造は類似する）の場合は、手動でテンプレートを作成し、或いは教師なし、教師ありの方法によりウェブページに含まれる製品情報の構造テンプレートを学習して、これらの学習により得られた構造テンプレートを用いて他の類似のウェブページを解析してもよい。図1Aに示すように、携帯電話のページの構造情報を学習することで、図書及び靴の製品情報を抽出してもよい。

10

【0004】

2種類目の方法は、単一の構造を有する（非類似の）ページの場合は、図1Bに示すように、ウェブページの構造を動的に解析し、キーワードのリストにより関連情報のウェブページにおける位置を特定し、値を抽出してもよい。

【発明の概要】

【発明が解決しようとする課題】

【0005】

以下は、本発明の態様を基本的に理解させるために、本発明の簡単な概要を説明する。なお、この簡単な概要は、本発明を網羅的な概要ではなく、本発明のポイント又は重要な部分を意図的に特定するものではなく、本発明の範囲を意図的に限定するものではなく、後述するより詳細的な説明の前文として、単なる概念を簡単な形で説明することを目的とする。

20

【0006】

本発明は、ウェブページから情報を抽出する方法、装置及び記憶媒体を提供する。

【課題を解決するための手段】

【0007】

本発明の1つの態様では、ウェブページから情報を抽出する方法であって、前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成するステップと、前記木構造におけるナビゲーションバーノードを決定するステップと、前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するステップと、マッチする葉ノードに対応するページにおける情報を抽出するステップと、を含む、方法を提供する。

30

【0008】

本発明のもう1つの態様では、ウェブページから情報を抽出する装置であって、前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成する木構造生成手段と、前記木構造におけるナビゲーションバーノードを決定するナビゲーションバーノード決定手段と、前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するマッチノード決定手段と、マッチする葉ノードに対応するページにおける情報を抽出する情報抽出手段と、を含む、装置を提供する。

40

【0009】

本発明の他の態様では、対応するコンピュータプログラムコード、コンピュータ読み取り可能な記憶媒体及びコンピュータプログラムプロダクトをさらに提供する。

【0010】

本発明に係るウェブページから情報を抽出する方法及び装置によれば、ホームページのURL（ユニフォームリソースロケータ）に基づいて、同一のドメイン名において分布している複数のウェブページから必要な情報を抽出することができる。

50

【 0 0 1 1 】

以下は図面を参照しながら本発明の好ましい実施形態を詳細に説明することにより、本発明の上記及び他の利点はより明確になる。

【 図面の簡単な説明 】

【 0 0 1 2 】

本開示の上記及び他の利点及び特徴を理解させるために、以下は図面を参照しながら本開示の具体的な実施形態を詳細に説明する。図面及び以下の詳細な説明は本明細書に含まれ、本明細書の一部を構成する。同一の機能及び構造を有する素子は同一の符号で示される。なお、これらの図面は単なる本開示の典型的な例を説明するためのものであり、本開示の範囲を限定するものではない。

【 図 1 A 】 類似の構造を有するウェブページの例を示す図である。

【 図 1 B 】 単一の構造を有するウェブページの情報抽出の例を示す図である。

【 図 2 A 】 複数のページの情報抽出の例を示す図である。

【 図 2 B 】 本発明の方法の全体的な流れの例を示す図である。

【 図 3 】 本発明の実施形態に係るウェブページから情報を抽出する方法の流れを示すフローチャートである。

【 図 4 A 】 ナビゲーションパーノードに対応する HTML 構造及び DOM 木構造の例を示す図である。

【 図 4 B 】 情報抽出の例を示す図である。

【 図 5 】 本発明の実施形態に係るウェブページから情報を抽出する装置の例を示すブロック図である。

【 図 6 】 本発明の実施形態に係る方法及び / 又は装置を実現可能な汎用パーソナルコンピュータの例示的な構成を示すブロック図である。

【 発明を実施するための形態 】

【 0 0 1 3 】

以下、図面を参照しながら本発明の例示的な実施例を詳細に説明する。説明の便宜上、明細書には実際の実施形態の全ての特徴が示されていない。なお、実際に実施する際に、開発者の具体的な目標を実現するために、特定の実施形態を変更してもよい、例えばシステム及び業務に関する制限条件に応じて実施形態を変更してもよい。また、開発作業が非常に複雑であり、且つ時間がかかるが、本公開の当業者にとって、この開発作業は単なる例の作業である。

【 0 0 1 4 】

なお、本発明を明確にするために、図面には本発明に密に関連する装置の構成要件及び / 又は処理のステップのみが示され、本発明と関係のない細部が省略されている。

【 0 0 1 5 】

上述したように、インターネットから情報を収集、抽出することは、知識ベースを構築する重要な手段である。図 1 A 及び図 1 B に示す従来技術は、一部の要求を満たすことができるが、依然として限界がある。

【 0 0 1 6 】

図 2 A に示すように、`http://owtware.com` は会社のホームページの URL であり、製品、協力会社、連絡先などの会社の情報は異なるページに分布し、3つのページの主要情報の所在する部分も類似の構造を有しない。

【 0 0 1 7 】

ホームページの URL のみが既知である場合、従来方法は、このような複数のページに分布している情報を抽出することができない。一方、通常ホームページの URL は容易に入手できる。このため、ホームページの URL 情報を拡張して他の情報を抽出する方法は、依然として解決すべき問題である。

【 0 0 1 8 】

従来技術に存在する問題を解決するために、本発明は、ホームページ URL のみが既知である場合、(1) 関連情報を含む他のページを自動的に拡張し、(2) 各関連ページが

10

20

30

40

50

ら主要情報を含む位置を取得し、(3)異なる属性タイプを有するページについて個別の情報抽出を行うことができる、複数のページに基づく情報抽出方法を提供する。

【0019】

図2Bは本発明の方法の全体的な流れの例を示す図である。図2Bに示すように、本発明に係る方法は、主に以下の3つの部分を含む。

【0020】

(1)ホームページを拡張することで複数のページの集合を取得する。

【0021】

(2)統計的方法を用いてウェブページの集合に対して統計的な分類を行い、ナビゲーションバーノード(navigation bar node)を取得し、そして、キーワード辞書を用いてナビゲーションバーノードに含まれる葉ノードのテキストのマッチングを行い、マッチするノード情報に基づいて抽出すべきページを取得する。

10

【0022】

(3)抽出すべきページの情報タイプに応じて、異なる解析器を用いて抽出を行う。

【0023】

以下は、図3、図4A及び図4Bを参照しながら、本発明の実施形態に係るウェブページから情報を抽出する方法を詳細に説明する。

【0024】

図3は本発明の実施形態に係るウェブページから情報を抽出する方法の流れを示すフローチャートである。

20

【0025】

まず、ステップ301において、ウェブページ及びその全ての拡張ウェブページにおける該ウェブページのドメイン名を含む各ページについて木構造を生成する。具体的には、本実施形態では、図2Aに示すURLを一例にすると、会社ホームページURLは $u_{root} = \text{http}://\text{www.owtware.com}/$ であり、抽出すべき情報は該会社の他の属性、例えば製品、連絡先などである。

【0026】

まず、クローラー(crawler)を用いて u_{root} に対応するHTMLページ p_{root} をクローリングし、ページを解析してそれに含まれる全てのURLの集合 $u = [u_0, u_1, u_2, \dots, u_n]$ を取得する。ページに含まれるURLが該会社に関連する場合があります、関連しない場合もあり、例えば広告や外部リンクなどの場合もあると考慮すると、特定のルールに従って一部のURLの集合 $u' = [u_0', u_1', u_2', \dots, u_n']$ を選択し、ここで、 u_i' には $\text{domain}(u_{root})$ が含まれ、 $\text{domain}(URL)$ はURLトップレベルドメイン名を抽出する操作であり、例えば $\text{domain}(u_{root}) = \text{www.owtware.com}$ 。このように、同一のドメイン名を有する全てのURL、例えば $\text{http}://\text{www.owtware.com}/\text{index.php}/\text{zh}/\text{products}/$ を保持することができる。

30

【0027】

好ましくは、 u_i' に対応するページ p_i は他のURL情報を含む可能性があると考えられるため、 p_i をさらに拡張してもよい。各 p_i について、同様のルールでURL及び対応するページを拡張し、毎回の拡張の後に同一のURL及びページを併合する。拡張のプロセスはn回だけ繰り返してもよい。一定の数のページを取得でき、且つページの数が多すぎないように、通常 $n = 2$ にしてもよい。これによって、同一のドメイン名を有するページの集合 $p = [\langle p_0, u_0 \rangle, \langle p_1, u_1 \rangle, \langle p_2, u_2 \rangle, \dots, \langle p_n, u_n \rangle]$ を取得でき、ここで、 p_i はウェブページを表し、 u_i はウェブページに対応するURLを表す。

40

【0028】

次に、ステップ302において、木構造におけるナビゲーションバーノードを決定する。具体的には、本実施形態では、集合 p からナビゲーションバーノードを取得する。上述したように、目的は、集合 p から該会社情報を含むページ、例えば製品、連絡先などを取

50

得することである。通常、ナビゲーションバーノードにおけるリンクにより、これらの情報に対応するページを取得できる。ナビゲーションバーノードを情報アンカーとして選択する主な理由は3つある。

【0029】

(1) 情報は正確である。ナビゲーションバーノードに含まれるリンクが指向するページは、会社の紹介と見なすことができる。例えば、「製品とサービス」に対応するページは該会社の製品を紹介し、「連絡先」は会社の住所、電話番号などの情報のページにリンクする。ウェブページにおける他の部分に出現するリンクは、必ずしも該会社の情報を説明するものではなく、他の会社の紹介や広告などの情報である可能性がある。

【0030】

(2) 情報は全面的である。ナビゲーションバーノードは基本的に該会社に関連する全ての情報を含み、ナビゲーションバーノードを取得すると、関連情報を含む全てのページを取得でき、これは後続の情報抽出に非常に役に立つ。

【0031】

(3) 比較的取得しやすい。異なるウェブページは異なる構造を有する可能性があるが、ナビゲーションバーノードの様式は殆ど同じである。このような共通性により、ウェブ構造からナビゲーションバーノードの位置を正確に見つけることができる。

【0032】

以下は、ナビゲーションバーノードの決定方法を例示的に説明する。

【0033】

上記の3つの特徴により、各ページ p_i ($p_i \in p$) におけるノードを計数することで、頻繁に出現するノードを取得してもよい。これらのノードにはナビゲーションバーノードが含まれるため、特徴値に基づいてこれらの頻繁に出現するノードを並び替えることでナビゲーションバーノードを取得してもよい。具体的な方法は以下の通りである。

【0034】

図4Aに示すように、集合 p における各ページ p_i について、まず p_i を Dom木の構造に変換する。

【0035】

Dom木における各葉ノード $node_i$ について、 $node_i$ の経路パターン $path_i$ を取得し、 $path_i$ は、該葉ノードに対応するテキストと、 n 番目の先祖ノードまでの経路により構成される。実際の経験によると、殆どのページでは、 n は5以上の整数値であってもよい。例えば、ナビゲーションバーノード「連絡先」について、 $n=5$ の場合は、 $path_i = \text{「ul_li_ul_li_a_連絡先」}$ を取得できる。

【0036】

次に、各 $path_i$ の文書頻度 df_i 、即ち $path_i$ が異なる文書に出現する回数を算出する。統計により経路頻度辞書 $node_pattern_dictionary \{ \langle path_1, df_1 \rangle, \dots, \langle path_n, df_n \rangle \}$ を取得してもよく、ここで、 $df_i > t$ 、 t は次のように設定された閾値である。

【数1】

$$t = \begin{cases} |p| * 0.2 & \text{if } |p| \geq 20 \\ |p| * 0.3 & \text{if } |p| \geq 10 \\ 3 & \text{else} \end{cases}$$

【0037】

ページ数 $|p|$ の最終結果への影響を低減するために、閾値 t を段階的に設定する。

【0038】

10

20

30

40

50

経路頻度辞書を取得した後、集合 p における各 p_i に対応する Dom 木構造に対して 2 回目の走査を行い、今回は、各非葉ノード $node_i$ について、それによりカバーされる全ての NULL でない葉ノードの集合が $c = [c_0, c_1, c_2, \dots, c_n]$ となると仮定すると、各 c_i について、 $path_i(c_i)$ が経路頻度辞書 $node_pattern_dictionary$ に存在する場合、該 $node_i$ の情報を記録する。最後に、候補辞書 $candidate_pattern_dictionary \{ \langle path_1, [df_1, cn_1] \rangle, \dots, \langle path_n, [df_n, cn_n] \rangle \}$ を取得してもよく、ここで、 $path_i$ は非葉ノード $node_i$ から先祖ノードまでの経路情報を表し、 df_i は文書頻度を表し、 cn_i は $node_i$ によりカバーされる全ての NULL でない葉ノードの数を表す。葉ノードの $path_i$ とは異なって、非葉ノードの $path_i$ はテキスト情報を含まない。図 4 A における 3 に示すように、「連絡先」から u_1 ノードまでの対応する経路は $u_1_li_u_1_div_div$ であり、 $n = 5$ となる。

【0039】

最後に、 $(cn * df / |p|)$ の値に従って候補辞書 $candidate_pattern_dictionary$ を並び替え、最大値に対応する経路をナビゲーションバーノード経路テンプレートとして取得し、該最大値に対応する経路における先祖ノードをナビゲーションバーノードとして決定してもよい。該会社のホームページの下にある所定の HTML ページについて、該テンプレートを用いてナビゲーションバーノードの位置を特定してもよい。

【0040】

なお、上記の統計的方法を用いてナビゲーションバーノードを決定することは、単なるナビゲーションバーノードの決定方法の一例である。本発明は、これに限定されず、他の適切な方法を用いてナビゲーションバーノードを決定してもよい。

【0041】

次に、ステップ 303 において、ナビゲーションバーノードによりカバーされる、1 つ又は複数のキーワードにマッチする葉ノードを決定する。具体的には、本実施形態では、ステップ 302 においてナビゲーションバーノードが取得された後に、該ナビゲーションバーノードによりカバーされる各 NULL でない葉ノードについて、辞書 $keyword_dict$ を用いて葉ノードに対応するテキストのマッチングを行う。辞書 $keyword_dict$ には、例えば「製品紹介」、「連絡先」などの所定のキーワードが含まれる。葉ノードがキーワードにマッチする場合、対応する HTML 要素から「href」属性を検索してもよく、その属性値は対応するウェブページの URL である。例えば、図 4 A における「連絡先」ノードに対応する HTML 要素には次のリンクが含まれる。

【0042】

`href = http://www.owtware.com/index.php/zh/about/contact-us/`

従って、集合 p から関連情報を含むウェブページの集合 $p' = [\langle p'_0, u'_0, t'_0 \rangle, \langle p'_1, u'_1, t'_1 \rangle, \langle p'_2, u'_2, t'_2 \rangle, \dots, \langle p'_n, u'_n, t'_n \rangle]$ を選択してもよく、ここで、 p'_i 及び u'_i は上記の定義された p_i 及び u_i と同じであり、 t'_i は、該ページに対応するタイプ、例えば製品、人物、連絡先などを表す。これによって、ページの異なるタイプに応じて、異なる解析器を選択して抽出を行うことができる。

【0043】

各 p'_i について、まず、HTML ページを前処理する必要がある。前処理の目的は、まずページにおける主要情報を抽出することである。このプロセスは共通のものであり、ウェブページのタイプ t' とは関係がない。抽出された結果は、後で抽出を行う時の入力としてもよい。図 4 B の (1) に示すように、元の HTML ページには多くの内容が含まれているが、実線の枠で示される部分のみが必要な内容であり、ナビゲーションバーノード、サイドリスト、ラベル Footer などの要素を含む他の部分を全て除去する必要がある。除去しないと、抽出時にノイズデータの影響を受けやすくなる。

【0044】

10

20

30

40

50

ステップ302において生成された経路頻度辞書 `node_pattern_dictionary` 及び候補辞書 `candidate_pattern_dictionary` を考慮すると、以下の方法を用いてナビゲーションバーノードによりカバーされる1つ又は複数のキーワードにマッチする葉ノードを決定してもよい。

【0045】

集合 p_i における非葉ノード $node_i$ について、それによりカバーされる全ての `NULL` でない葉ノードの集合が $c = [c_0, c_1, c_2, \dots, c_n]$ であると仮定すると、次の3つの条件が同時に満たされた場合、 $node_i$ が1つ又は複数のキーワードにマッチする葉ノードを含むターゲット内容ノードであると決定してもよい。

【数2】

- ✓ $path(node_i) \notin candidate_pattern_dictionary,$
- ✓ $path(c_i) \notin node_pattern_dictionary, \forall c_i \in c, \text{ 且つ}$
- ✓ $\sum text_len(c_i) > \sum text_len(c_j)$

【0046】

ここで、 c_i は $node_i$ によりカバーされる `NULL` でない葉ノードであり、 c_j は $node_j$ によりカバーされる `NULL` でない葉ノードであり、 $i \neq j$ となり、 $text_len(*)$ は葉ノードに対応するテキストの長さを表す。言い換えれば、 $node_i$ によりカバーされる全ての `NULL` でない葉ノードのテキストの合計長さは、他のノード $node_j$ によりカバーされる全ての `NULL` でない葉ノードのテキストの合計長さよりも大きい。

【0047】

上記の3つの条件を同時に満たすノード $node_i$ が決定されると、所定のキーワードにマッチする葉ノードが決定されることを意味する。

【0048】

最後に、ステップ304において、マッチする葉ノードに対応するページにおける情報を抽出する。具体的には、本実施形態では、上記3つの条件を満たすノード $node_i$ が決定された後、該ノードによりカバーされる葉ノードに含まれる情報を抽出してもよい。

【0049】

好ましくは、その各葉ノードを独立した属性抽出空間としてもよく、図4Bにおける(2)及び(3)に示すように、各ノード `<div class="panel-grid-cell" ... >` を独立した属性空間とする。これによって、属性値の境界を決定することができ、即ち、各値はセクション `{{...}}` からの値のみである。例えば、人物情報を抽出する場合、セクション `{{...}}` に含まれる情報は同一の人物を表すためのものであり、異なる `{{...}}` の情報は異なる人物を表すと見なしてもよい。そのため、抽出エラーを回避することができる。

【0050】

好ましくは、抽出範囲が決定された後、 p_i のタイプ t_i に応じて、異なる解析器、例えばエンティティ認識器 (NER)、固有名詞認識器、数値認識器などを選択して特定情報の抽出を行ってもよい。図4Bの(3)では、固有名詞認識器の結果の例を示している。

【0051】

なお、以上は会社ホームページに基づいて関連情報を抽出することを説明しているが、本発明はこれに限定されず、必要に応じて任意のウェブページの任意の情報の抽出に拡張されてもよい。

【0052】

上記の方法は、コンピュータ実行可能なプログラムにより完全に実現されてもよいし、ハードウェア及び/又はファームウェアを用いて部分的又は完全に実現されてもよい。ハ

10

20

30

40

50

ードウェア及び/又はファームウェアにより実現される場合、又はコンピュータ実行可能なプログラムがプログラムを実行可能なハードウェア装置にロードされる場合、後述するウェブページから情報を抽出する装置が実現される。以下は、上述した詳細な内容を省略し、これらの装置の概要を説明する。なお、これらの装置は上記の方法を実行することができるが、上記方法は後述する装置の構成部を採用し、或いは構成部により実行されるものに限定されない。

【0053】

図5は本発明の実施形態に係るウェブページから情報を抽出する装置500の例を示すブロック図である。装置500は、木構造生成部501、ナビゲーションパーノード決定部502、マッチノード決定部503及び情報抽出部504を含む。木構造生成部501は、ウェブページ及びその全ての拡張ウェブページにおける該ウェブページのドメイン名を含む各ページについて木構造を生成する。ナビゲーションパーノード決定部502は、該木構造におけるナビゲーションパーノードを決定する。マッチノード決定部503は、該ナビゲーションパーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定する。情報抽出部504は、マッチする葉ノードに対応するページにおける情報を抽出する。

10

【0054】

図5に示すウェブページから情報を抽出する装置500は図3に示す方法に対応する。よって、ウェブページから情報を抽出する装置500の詳細は、図3におけるウェブページから情報を抽出する方法について説明において既に詳細に説明され、ここでその説明を省略する。

20

【0055】

上記処理及び装置はソフトウェア及び/又はファームウェアにより実現されてもよい。ソフトウェア及び/又はファームウェアにより実施されている場合、記憶媒体又はネットワークから専用のハードウェア構成を有するコンピュータ(例えば図6示されている汎用パーソナルコンピュータ600)に上記方法を実施するためのソフトウェアを構成するプログラムをインストールしてもよく、該コンピュータは各種のプログラムがインストールされている場合は各種の機能などを実行できる。

【0056】

図6は本発明の実施形態に係る方法及び/又は装置を実現可能な汎用パーソナルコンピュータの例示的な構成を示すブロック図である。図6において、中央処理部(CPU)601は、読み出し専用メモリ(ROM)602に記憶されているプログラム、又は記憶部608からランダムアクセスメモリ(RAM)603にロードされたプログラムにより各種の処理を実行する。RAM603には、必要に応じて、CPU601が各種の処理を実行するに必要なデータが記憶されている。CPU601、ROM602、及びRAM603は、バス604を介して互いに接続されている。入力/出力インターフェース605もバス604に接続されている。

30

【0057】

入力部606(キーボード、マウスなどを含む)、出力部607(ディスプレイ、例えばブラウン管(CRT)、液晶ディスプレイ(LCD)など、及びスピーカなどを含む)、記憶部608(例えばハードディスクなどを含む)、通信部609(例えばネットワークのインタフェースカード、例えばLANカード、モデムなどを含む)は、入力/出力インターフェース605に接続されている。通信部609は、ネットワーク、例えばインターネットを介して通信処理を実行する。必要に応じて、ドライバ610は、入力/出力インターフェース605に接続されてもよい。取り外し可能な媒体611は、例えば磁気ディスク、光ディスク、光磁気ディスク、半導体メモリなどであり、必要に応じてドライバ610にセットアップされて、その中から読みだされたコンピュータプログラムは必要に応じて記憶部608にインストールされている。

40

【0058】

ソフトウェアにより上記処理を実施する場合、ネットワーク、例えばインターネット、

50

又は記憶媒体、例えば取り外し可能な媒体 6 1 1 を介してソフトウェアを構成するプログラムをインストールする。

【 0 0 5 9 】

なお、これらの記憶媒体は、図 6 に示されている、プログラムを記憶し、機器と分離してユーザへプログラムを提供する取り外し可能な媒体 6 1 1 に限定されない。取り外し可能な媒体 6 1 1 は、例えば磁気ディスク（フロッピーディスク（登録商標）を含む）、光ディスク（光ディスク - 読み出し専用メモリ（CD-ROM）、及びデジタル多目的ディスク（DVD）を含む）、光磁気ディスク（ミニディスク（MD）（登録商標））及び半導体メモリを含む。或いは、記憶媒体は、ROM 6 0 2、記憶部 6 0 8 に含まれるハードディスクなどであってもよく、プログラムを記憶し、それらを含む機器と共にユーザへ提供される。

10

【 0 0 6 0 】

本発明は、対応するコンピュータプログラムコード、機器が読み取り可能な命令コードが記憶されているコンピュータプログラムプロダクトをさらに提供する。該命令コードは、機器により読み取られ、実行される際に、上記の本発明の実施例に係る方法を実行することができる。

【 0 0 6 1 】

それに応じて、本発明は、機器が読み取り可能な命令コードを含むプログラムプロダクトが記録されている記憶媒体をさらに含む。該記憶媒体は、フロッピーディスク、光ディスク、光磁気ディスク、メモリカード、メモリスティック等を含むが、これらに限定されない。

20

【 0 0 6 2 】

また、上述の各実施例を含む実施形態に関し、更に以下の付記を開示する。

（付記 1）

ウェブページから情報を抽出する方法であって、

前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成するステップと、

前記木構造におけるナビゲーションバーノードを決定するステップと、

前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するステップと、

30

マッチする葉ノードに対応するページにおける情報を抽出するステップと、を含む、方法。

（付記 2）

統計的方法を用いて前記ナビゲーションバーノードを決定する、付記 1 に記載の方法。

（付記 3）

前記木構造におけるナビゲーションバーノードを決定するステップは、

前記木構造に出現する回数が所定閾値よりも大きい葉ノードのみを含む非葉ノードを決定するステップと、

前記非葉ノードを並び替えて前記ナビゲーションバーノードを決定するステップと、を含む、付記 2 に記載の方法。

40

（付記 4）

葉ノードの出現回数が所定閾値よりも大きいか否かを決定することは、

前記葉ノードのテキスト及び経路情報の前記木構造における出現回数が前記所定閾値よりも大きいか否かを決定すること、を含む、付記 3 に記載の方法。

（付記 5）

前記経路情報は、前記葉ノードからその n 番目の先祖ノードまでの経路であり、n は正整数である、付記 4 に記載の方法。

（付記 6）

n は 5 以上である、付記 5 に記載の方法。

（付記 7）

50

前記非葉ノードを並び替えて前記ナビゲーションバーノードを決定するステップは、
前記非葉ノードの特徴値を計算するステップであって、前記特徴値は、前記非葉ノードによりカバーされる葉ノードの数及び前記回数により決定される、ステップと、
前記非葉ノードのうちの最大の特徴値を有する非葉ノードを前記ナビゲーションバーノードとして決定するステップと、を含む、付記 3 に記載の方法。

(付記 8)

前記特徴値は、前記非葉ノードによりカバーされる葉ノードの数と前記回数との積の、前記ウェブページのドメイン名を含むページの総数に対する比率である、付記 7 に記載の方法。

(付記 9)

マッチする葉ノードに対応するページにおける情報を抽出するステップは、
前記マッチする葉ノードに対応するページに含まれるターゲットノードを決定するステップと、
前記ターゲットノードによりカバーされる各葉ノードのテキストをそれぞれ抽出するステップと、を含む、付記 1 乃至 8 の何れかに記載の方法。

(付記 10)

前記ターゲットノードは、
前記ターゲットノードに含まれる各葉ノードのテキスト及び経路情報の前記木構造における出現回数が前記所定閾値以下であること、
前記ターゲットノードが、前記木構造に出現する回数が所定閾値よりも大きい葉ノードのみを含む非葉ノードのうちの非葉ノードではないこと、及び

前記ターゲットノードに含まれる全ての葉ノードのテキストの合計長さが該木構造における他の非葉ノードのテキストの合計長さよりも大きいこと、により決定される、付記 9 に記載の方法。

(付記 11)

前記ターゲットノードによりカバーされる各葉ノードのテキストをそれぞれ抽出するステップは、
前記ターゲットノードに対応するページのタイプに応じて、異なる解析器を選択して抽出を行うステップ、を含む、付記 9 に記載の方法。

(付記 12)

前記ターゲットノードの各葉ノードを独立した属性抽出空間とする、付記 11 に記載の方法。

(付記 13)

前記解析器は、エンティティ認識器、固有名詞認識器又は数値認識器である、付記 11 に記載の方法。

(付記 14)

決定されたナビゲーションバーノードの経路情報を用いて前記ウェブページ及びその全ての拡張ウェブページにおけるナビゲーションバーノードを決定する、付記 1 乃至 8 の何れかに記載の方法。

(付記 15)

URL トップレベルドメイン名を抽出することにより、前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含むページを決定する、付記 1 乃至 8 の何れかに記載の方法。

(付記 16)

前記木構造は、HTML 文書オブジェクトモデル (DOM) である、付記 1 乃至 8 の何れかに記載の方法。

(付記 17)

前記キーワードは、所定のキーワードである、付記 1 乃至 8 の何れかに記載の方法。

(付記 18)

前記拡張ウェブページを n 回だけ拡張して前記ウェブページのドメイン名を含むページ

10

20

30

40

50

を取得し、nは2以上の整数である、付記1乃至8の何れかに記載の方法。

(付記19)

ウェブページから情報を抽出する装置であって、

前記ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成する木構造生成手段と、

前記木構造におけるナビゲーションバーノードを決定するナビゲーションバーノード決定手段と、

前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するマッチノード決定手段と、

マッチする葉ノードに対応するページにおける情報を抽出する情報抽出手段と、を含む装置。

10

(付記20)

プログラムを記憶したコンピュータ読み取り可能な記憶媒体であって、前記プログラムがプロセッサにより実行される際に、

ウェブページ及びその全ての拡張ウェブページにおける前記ウェブページのドメイン名を含む各ページについて木構造を生成するステップと、

前記木構造におけるナビゲーションバーノードを決定するステップと、

前記ナビゲーションバーノードによりカバーされる、1つ又は複数のキーワードにマッチする葉ノードを決定するステップと、

マッチする葉ノードに対応するページにおける情報を抽出するステップと、を実行させる、記憶媒体。

20

【0063】

なお、用語「含む」、「有する」又は他の任意の変形は、排他的に含むことに限定されず、一連の要素を含むプロセス、方法、物又は装置は、これらの要素を含むことだけではなく、明示的に列挙されていない他の要素、又はこのプロセス、方法、物若しくは装置の固有の要素を含む。また、さらなる制限がない限り、用語「1つの...を含む」より限定された要素は、該要素を含むプロセス、方法、物又は装置に他の同一の要素が存在することを排除しない。

【0064】

以上は図面を参照しながら本発明の好ましい実施例を説明しているが、上記実施例及び例は例示的なものであり、制限的なものではない。当業者は、特許請求の範囲の主旨及び範囲内で本発明に対して各種の修正、改良、均等的なものに変更してもよい。これらの修正、改良又は均等的なものに変更することは本発明の保護範囲に含まれるものである。

30

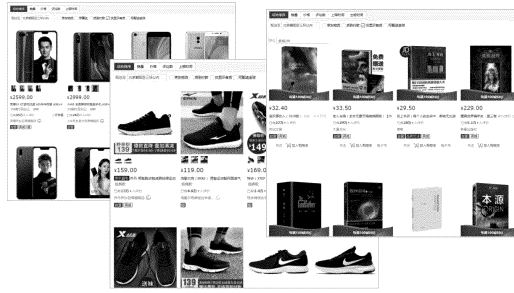
40

50

【図面】

【図 1 A】

類似の構造を有するウェブページの例を示す図



【図 1 B】

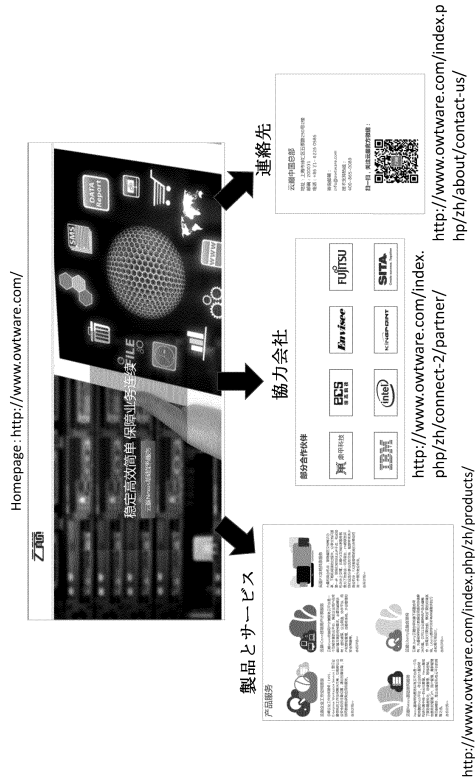
単一の構造を有するウェブページの情報抽出の例を示す図

キーワード	値
社名	富士通株式会社 (FUJITSU LIMITED) [法人番号: 1020000744012]
所在地	本店 本社事務所 住所: 〒211-8501 神奈川県川崎市 中原区上田4-1-1 富士通ビルディング Tel: 044-777-1111 (国内) Tel: 03-652-2220 (海外)
代表者	代表取締役社長 日中通信(たけ) 隆雄
設立	1959年6月2日 [富士通(株)創業]
事業内容	ソフトウェアソリューション ハードウェアソリューション サービスソリューション

10

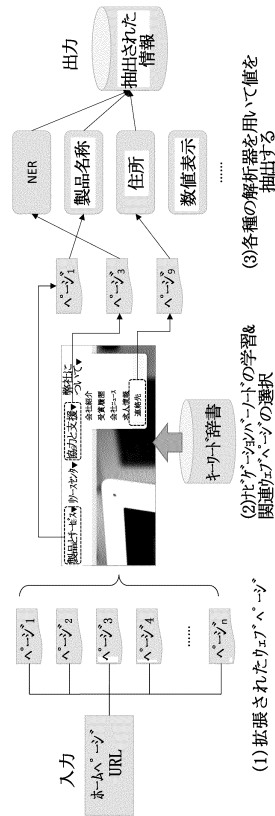
【図 2 A】

複数のページの情報の抽出の例を示す図



【図 2 B】

本発明の方法の全体的な流れの例を示す図



20

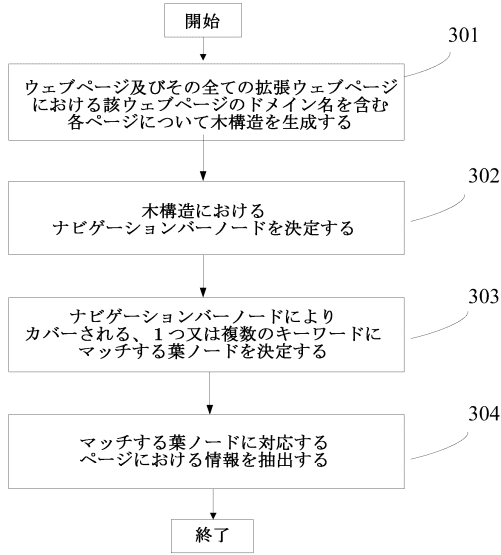
30

40

50

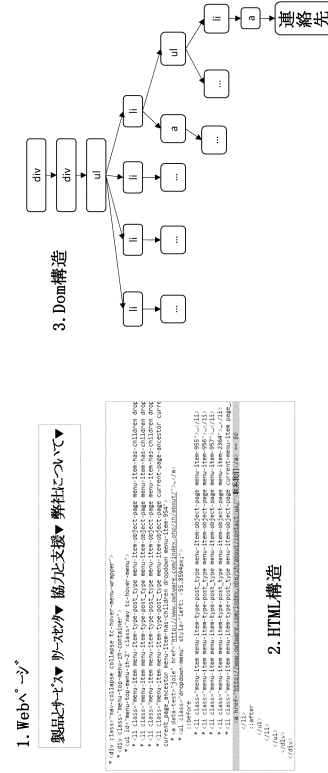
【 図 3 】

本発明の実施形態に係るウェブページから情報を抽出する方法の流れを示すフローチャート



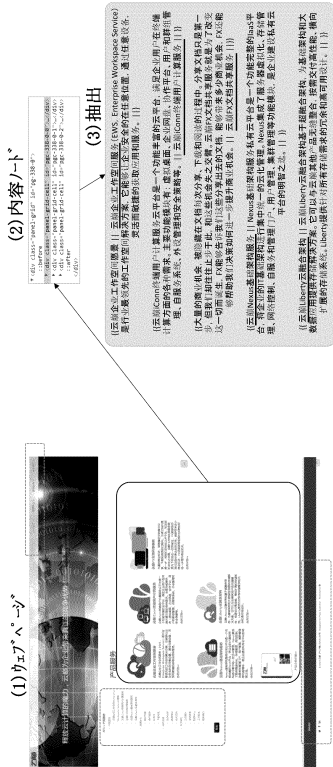
【 図 4 A 】

ナビゲーションパーノードに対応するHTML構造及びDOM木構造の例を示す図



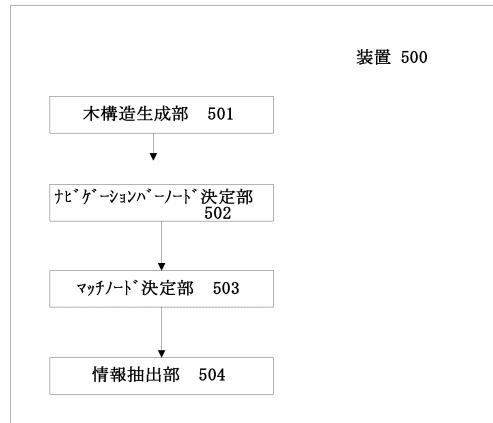
【 図 4 B 】

情報抽出の例を示す図



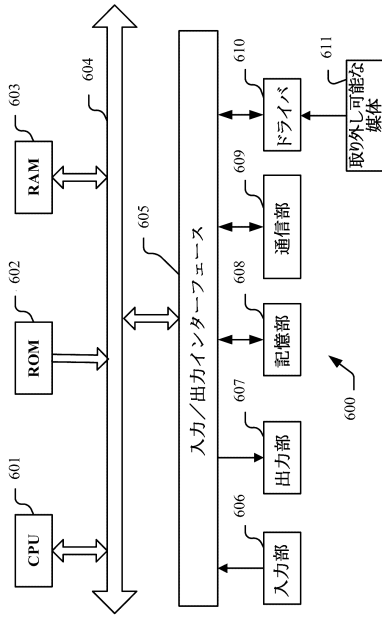
【 図 5 】

本発明の実施形態に係るウェブページから情報を抽出する装置の例を示すブロック図



【図6】

本発明の実施形態に係る方法及び/又は装置を実現可能な汎用パーソナルコンピュータの例示的な構成を示すブロック図



10

20

30

40

50

フロントページの続き

(72)発明者 孟 遥
中国, 100027, ベイジン, チャオヤン ディストリクト, ゴン ティ ベイ ルウ ナンバー 2
エイ, パシフィック センチュリー プレイス, スペース 8, ゲート 6, ユニット 3エフ, 35
5 富士通研究開発中心有限公司内

(72)発明者 孫 俊
中国, 100027, ベイジン, チャオヤン ディストリクト, ゴン ティ ベイ ルウ ナンバー 2
エイ, パシフィック センチュリー プレイス, スペース 8, ゲート 6, ユニット 3エフ, 35
5 富士通研究開発中心有限公司内

審査官 甲斐 哲雄

(56)参考文献 特開2009-042908(JP, A)
特開2016-201112(JP, A)
中国特許出願公開第105069107(CN, A)
中国特許出願公開第103823824(CN, A)

(58)調査した分野 (Int.Cl., DB名)
G06F 16/00 - 16/958
H04L 67/02