

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6834774号  
(P6834774)

(45) 発行日 令和3年2月24日(2021.2.24)

(24) 登録日 令和3年2月8日(2021.2.8)

(51) Int. Cl. F 1  
G 0 6 F 16/29 (2019.01) G 0 6 F 16/29

請求項の数 1 (全 11 頁)

|   |  |
|---|--|
| <p>(21) 出願番号 特願2017-101200 (P2017-101200)<br/>                 (22) 出願日 平成29年5月22日 (2017.5.22)<br/>                 (65) 公開番号 特開2018-195272 (P2018-195272A)<br/>                 (43) 公開日 平成30年12月6日 (2018.12.6)<br/>                 審査請求日 令和1年8月23日 (2019.8.23)</p> | <p>(73) 特許権者 000003207<br/>                 トヨタ自動車株式会社<br/>                 愛知県豊田市トヨタ町1番地<br/>                 (74) 代理人 100107766<br/>                 弁理士 伊東 忠重<br/>                 (74) 代理人 100070150<br/>                 弁理士 伊東 忠彦<br/>                 (72) 発明者 鈴木 功一<br/>                 愛知県豊田市トヨタ町1番地 トヨタ自動車株式会社内<br/>                 審査官 ▲はま▼中 信行</p> |
|---|--|

最終頁に続く

(54) 【発明の名称】 情報抽出装置

(57) 【特許請求の範囲】

【請求項1】

文書データを取得する文書データ取得部と、  
 前記文書データから住所の候補文字列を抽出する候補文字列抽出部と、  
 前記候補文字列に対してジオコーディングを行うことにより、位置情報の取得を試みる位置情報取得部と、  
 前記位置情報取得部による前記位置情報の取得結果に応じて、前記住所とする文字列を決定する住所文字列決定部と  
 を備え、  
 前記位置情報取得部による前記位置情報の取得に成功した場合、前記位置情報取得部による前記位置情報の再取得に失敗するまで、  
 前記候補文字列抽出部が、前記候補文字列を後方に延長して、前記文書データから前記候補文字列を再抽出し、  
 前記位置情報取得部が、再抽出された前記候補文字列に対して前記ジオコーディングを行うことにより、前記位置情報の再取得を試み、  
 前記位置情報取得部による前記位置情報の再取得に失敗した場合、  
 前記住所文字列決定部が、前記位置情報の再取得に失敗する直前に、前記位置情報の取得に成功した前記候補文字列を、前記住所とする文字列に決定する  
 情報抽出装置。

【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、情報抽出装置に関する。

## 【背景技術】

## 【0002】

従来より、WEBクローリングに代表されるように、WEBページ等の文書データから施設情報（例えば、POI（Point of Interest）名称、住所、電話番号、郵便番号、キーワード等）を抽出し、当該施設情報をデータベースに自動的に蓄積できるようにした技術が知られている。

## 【0003】

下記特許文献1には、インターネットに接続されたサーバから文書データを取得し、当該文書データに含まれている住所文字列に、「都道府県」、「市町村」、「町域」、「街区」、「号」のいずれまでが含まれているかによって、当該住所文字列の詳細度（1～5）を決定する技術が開示されている。例えば、住所文字列に「号」までが含まれている場合には、最も高い詳細度「5」が決定されるといった具合である。また、下記特許文献1には、所定の閾値以上の詳細度を有する住所文字列を含む文書データを、データベースに組み込むようにした技術が開示されている。

10

## 【先行技術文献】

## 【特許文献】

## 【0004】

【特許文献1】特開2012-256356号公報

20

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0005】

しかしながら、従来の技術では、文書データから抽出する住所文字列の終端を高精度に特定することができない。このため、従来の技術では、住所の途中までしか住所文字列として取得しない場合や、住所に続く住所以外の語句までも住所文字列として取得してしまう場合がある。また、従来の技術では、文書データから複数のPOI名称の候補文字列が抽出された場合、いずれの候補文字列が実際のPOI名称であるかを判断することが困難であるため、実際にはPOI名称ではない候補文字列を、POI名称として誤って抽出してしまう場合がある。

30

## 【0006】

このようなことから、従来、文書データからの施設情報の抽出精度を高めることが困難であった。

## 【0007】

本発明は、上述した従来技術の課題を解決するため、文書データからの施設情報の抽出処理を適切に行い、文書データからの施設情報の抽出精度を高めることを目的とする。

## 【課題を解決するための手段】

## 【0008】

本発明の実施形態の情報抽出装置は、文書データを取得する文書データ取得部と、前記文書データから住所の候補文字列を抽出する候補文字列抽出部と、前記候補文字列に対してジオコーディングを行うことにより、位置情報の取得を試みる位置情報取得部と、前記位置情報取得部による前記位置情報の取得結果に応じて、前記住所とする文字列を決定する住所文字列決定部とを備え、前記位置情報取得部による前記位置情報の取得に成功した場合、前記位置情報取得部による前記位置情報の再取得に失敗するまで、前記候補文字列抽出部が、前記候補文字列を後方に延長して、前記文書データから前記候補文字列を再抽出し、前記位置情報取得部が、再抽出された前記候補文字列に対して前記ジオコーディングを行うことにより、前記位置情報の再取得を試み、前記位置情報取得部による前記位置情報の再取得に失敗した場合、前記住所文字列決定部が、前記位置情報の再取得に失敗する直前に、前記位置情報の取得に成功した前記候補文字列を、前記住所とする文字列に決

40

50

定する。

【発明の効果】

【0009】

文書データからの施設情報の抽出処理を適切に行い、文書データからの施設情報の抽出精度を高めることができる。

【図面の簡単な説明】

【0010】

【図1】実施形態に係る情報抽出装置の機能構成を示す図である。

【図2】実施形態に係る情報抽出装置による処理の手順を示すフローチャートである。

【図3】実施形態に係るアノテータ処理部によるアノテータ処理の手順を示すフローチャートである。

10

【図4】実施形態に係るアノテータ処理部による住所取得処理の手順を示すフローチャートである。

【図5】実施形態に係るアノテータ処理部によるナイーブベイズ推定値取得処理の手順を示すフローチャートである。

【発明を実施するための形態】

【0011】

以下、図面を参照して、本発明の実施形態の情報抽出装置について説明する。

【0012】

(情報抽出装置100の機能構成)

20

図1は、実施形態に係る情報抽出装置100の機能構成を示す図である。図1に示す情報抽出装置100は、WEBクローリングを行うことによって、インターネット上のWEBページ110(「文書データ」の一例)からPOIデータ(「施設情報」の一例)を抽出し、当該POIデータを施設情報DB120へ登録および更新することが可能な装置である。

【0013】

図1に示すように、情報抽出装置100は、クローラ処理部101、パーサ処理部102、スクレーパ処理部103、アノテータ処理部104、およびデータリンカ処理部105を備える。

【0014】

30

クローラ処理部101は、WEBクローリングを行うことにより、インターネット上のWEBサイトからWEBページ110を取得し、当該WEBページ110をメモリに格納する。すなわち、クローラ処理部101は、「文書データ取得部」としての機能を有する。

【0015】

パーサ処理部102は、クローラ処理部101によって取得されたWEBページ110に対してパーサ処理を行うことにより、当該WEBページ110から、特定のキーワードによる、特定のHTML(HyperText Markup Language)ファイルの選択を行う。

【0016】

スクレーパ処理部103は、パーサ処理部102によって選択されたHTMLファイルに対してスクレーパ処理を行うことにより、当該HTMLファイルから不要部分を削除し、残りの部分を構造体として出力する。

40

【0017】

アノテータ処理部104は、スクレーパ処理部103によって出力された構造体に対してアノテータ処理を行うことにより、当該構造体に含まれるテキストデータを解析し、当該テキストデータから、予め定義された属性値(POI名称、住所、電話番号、郵便番号、キーワード等)を取得する。

【0018】

特に、アノテータ処理部104は、「候補文字列抽出部」、「位置情報取得部」、および「住所文字列決定部」としての機能を有している。すなわち、アノテータ処理部104

50

は、WEBページ110から住所の候補文字列を抽出し、当該候補文字列に対してジオコーディングを行うことにより位置情報の取得を試み、位置情報の取得結果に応じて、住所とする文字列を決定することができる。これにより、アナテータ処理部104は、WEBページ110から抽出する住所文字列の終端を高精度に特定することができる。この点については、図4を用いて詳細に説明する。

#### 【0019】

さらに、アナテータ処理部104は、WEBページ110からPOI名称の候補文字列を抽出し、各候補文字列について、ナイーブベイズ推定値を取得することができる。そして、アナテータ処理部104は、POI名称の候補文字列が、所定文字列を含む、または、強調されている場合、その候補文字列のナイーブベイズ推定値を高めることができる。これにより、アナテータ処理部104は、実際にPOI名称である可能性が最も高い文字列を、POI名称として抽出することができる。この点については、図5を用いて詳細に説明する。

10

#### 【0020】

データリンク処理部105は、データリンク処理を行うことにより、アナテータ処理部104により取得された各属性値をPOIデータとして、当該POIデータに対して、ジオコーディング、POIマスタとの名寄せ等を行い、当該POIデータを施設情報DB120に対して登録または更新する。

#### 【0021】

なお、情報抽出装置100の各機能は、例えば、各種情報処理装置（例えば、サーバ、パーソナルコンピュータ等）において、各種記憶装置（例えば、ROM（Read Only Memory）、フラッシュメモリ等）に記憶されたプログラムを、コンピュータ（例えば、CPU（Central Processing Unit）等）が実行することにより、実現される。

20

#### 【0022】

（情報抽出装置100による処理の手順）

図2は、実施形態に係る情報抽出装置100による処理の手順を示すフローチャートである。図2の処理は、例えば、情報抽出装置100にスケジュール設定されることにより、情報抽出装置100によって定期的（例えば、1日毎）に実行される。

#### 【0023】

まず、クローラ処理部101が、WEBクローリングを行うことにより、インターネット上のWEBサイトからWEBページ110を取得し、当該WEBページ110をメモリに格納する（ステップS201）。

30

#### 【0024】

次に、パーサ処理部102が、ステップS201で取得されたWEBページ110に対してパーサ処理を行うことにより、当該WEBページ110から、特定のキーワードによる、特定のHTMLファイルの選択を行う（ステップS202）。

#### 【0025】

次に、スクレーパ処理部103が、ステップS202で選択されたHTMLファイルに対してスクレーパ処理を行うことにより、当該HTMLファイルから不要部分を削除し、残りの部分（すなわち、POIデータの抽出対象とするテキストデータ。例えば、口コミ情報等）を構造体として出力する（ステップS203）。

40

#### 【0026】

次に、アナテータ処理部104が、ステップS203で出力された構造体に対してアナテータ処理を行うことにより、当該構造体に含まれるテキストデータを解析し、当該テキストデータから、予め定義された属性値（POI名称、住所、電話番号、郵便番号、キーワード等）を取得する（ステップS204）。なお、アナテータ処理部104によるアナテータ処理の詳細については、図3を用いて後述する。

#### 【0027】

次に、データリンク処理部105が、データリンク処理を行うことにより、ステップS204で取得された各属性値をPOIデータとして、当該POIデータに対して、ジオコ

50

ーディング、POIマスタとの名寄せ等を行い、当該POIデータを施設情報DBに対して登録または更新する(ステップS205)。そして、情報抽出装置100は、図2に示す一連の処理を終了する。

#### 【0028】

(アノテータ処理部104によるアノテータ処理の手順)

図3は、実施形態に係るアノテータ処理部104によるアノテータ処理の手順を示すフローチャートである。図3は、図2にフローチャートにおけるステップS204のアノテータ処理を詳細に説明するものである。図3の処理には、スクレーパ処理部103から出力された構造体(HTML構造)が入力される。

#### 【0029】

まず、アノテータ処理部104は、構造体に含まれるテキストデータに対して、形態素解析を行い、当該テキストデータを、複数の形態素(単語、品詞等)単位に分割する(ステップS301)。

#### 【0030】

次に、アノテータ処理部104は、ステップS301で複数の形態素に分割されたテキストデータの中から、郵便番号および電話番号を取得する(ステップS302)。例えば、アノテータ処理部104は、「」および数字からなる所定のフォーマットの文字列(例えば、「xxx-xxxx」)や、直前に「郵便番号」、「〒」等が存在する文字列を、郵便番号として取得する。また、例えば、アノテータ処理部104は、「」および数字からなる所定のフォーマットの文字列(例えば、「xxx-xxxx-xxxx」)や、直前に「電話番号」、「TEL」等が存在する文字列を、電話番号として取得する。

#### 【0031】

次に、アノテータ処理部104は、住所取得処理を実行することにより、ステップS301で複数の形態素に分割されたテキストデータの中から、住所を取得する(ステップS303)。住所取得処理の詳細については、図4を用いて後述する。

#### 【0032】

次に、アノテータ処理部104は、ステップS301で複数の形態素に分割されたテキストデータのうち、<title>タグが付されている部分と、<h>タグが付されている部分とのそれぞれに対して、POI名称の候補文字列の抽出を試みる(ステップS304)。<title>タグおよび<h>タグは、POI名称が設定されている可能性が高いからである。但し、これに限らず、アノテータ処理部104は、これ以外のタグが付されている部分についても、POI名称の候補文字列の抽出を試みるようにしてもよい。

#### 【0033】

次に、アノテータ処理部104は、ステップS304で抽出されたPOI名称の候補文字列から、POI名称として不要と思われる部分を除去する(ステップS305)。さらに、アノテータ処理部104は、頻出語フィルタ処理を行うことにより、ステップS304で抽出されたPOI名称の候補文字列の中から、POI名称である可能性の高い候補文字列を抽出する(ステップS306)。例えば、アノテータ処理部104は、POI名称の候補文字列に、予め学習しておいたPOI名称の頻出語が含まれている場合、その候補文字列がPOI名称である可能性が高いと判断する。

#### 【0034】

次に、アノテータ処理部104は、ナイーブベイズ推定値取得処理を実行することにより、ステップS306で抽出された各候補文字列に対して、ナイーブベイズ推定値を取得する(ステップS307)。ナイーブベイズ推定値取得処理の詳細については、図4を用いて後述する。そして、アノテータ処理部104は、図3に示す一連の処理を終了する。

#### 【0035】

(アノテータ処理部104による住所取得処理の手順)

図4は、実施形態に係るアノテータ処理部104による住所取得処理の手順を示すフローチャートである。図4は、図3にフローチャートにおけるステップS303の住所取得処理を詳細に説明するものである。なお、図4の処理には、複数のテキストデータを含む

10

20

30

40

50

構造体が入力される。これに応じて、アノテータ処理部104は、図4の処理を、構造体に含まれるテキストデータ毎に実行する。

【0036】

まず、アノテータ処理部104は、構造体に含まれるテキストデータから、住所を含むと推定される文書を抽出する(ステップS401)。例えば、アノテータ処理部104は、住所に関する特定のキーワード(例えば、「住所：」、「県」、「市」等)を含む文書(例えば、「この度ついに××県××市××〇丁目〇番地〇号に新規オープンしました」等)を抽出する。

【0037】

次に、アノテータ処理部104は、ステップS401で抽出された文書が、住所の書式(例えば、「××県××市」等)に合致する文字列を含んでいるか否かを判断する(ステップS402)。

【0038】

ステップS402において、住所の書式に合致する文字列を含んでいないと判断された場合(ステップS402:No)、アノテータ処理部104は、図4に示す一連の処理を終了する。

【0039】

一方、ステップS402において、住所の書式に合致する文字列を含んでいると判断された場合(ステップS402:Yes)、アノテータ処理部104は、住所の書式に合致すると判断された文字列の長さが、128文字未満であるか否かを判断する(ステップS403)。

【0040】

ステップS403において、住所の書式に合致すると判断された文字列の長さが、128文字未満ではないと判断された場合(ステップS403:No)、アノテータ処理部104は、図4に示す一連の処理を終了する。

【0041】

一方、ステップS403において、住所の書式に合致すると判断された文字列の長さが、128文字未満であると判断された場合(ステップS403:Yes)、アノテータ処理部104は、住所の書式に合致すると判断された文字列を候補文字列とし、当該候補文字列に対してジオコーディングを実施する(ステップS404)。例えば、アノテータ処理部104は、特定の機関から提供されたジオコーディング用のAPI(Application Programming Interface)を実行することにより、候補文字列に対応する位置情報(経度および緯度)を取得する。

【0042】

そして、アノテータ処理部104は、ジオコーディングによる位置情報の取得に成功したか否かを判断する(ステップS405)。ステップS405において、ジオコーディングによる位置情報の取得に成功したと判断された場合(ステップS405:Yes)、アノテータ処理部104は、位置情報の取得に成功した候補文字列をメモリに格納し(ステップS406)、候補文字列を後方に延長して(ステップS407)、ステップS401で抽出された文書から、候補文字列を再抽出する(ステップS408)。そして、アノテータ処理部104は、ステップS404に処理を戻す。

【0043】

なお、アノテータ処理部104は、例えば、ステップS407による候補文字列の延長を、住所に関する特定の語句単位(例えば、都道府県、市区町村、番地等)で行う。例えば、現在の候補文字列が「××県」であった場合において、その次の語句が「××市」であった場合、アノテータ処理部104は、「××県××市」を新たな候補文字列としてもよい。ここで、アノテータ処理部104は、現在の候補文字列の次に、住所に関する特定の語句ではない品詞が存在する場合、候補文字列を、その品詞まで延長してもよい。例えば、現在の候補文字列が「××県××市××〇丁目〇番地〇号」であった場合において、その次の品詞が「に」であった場合、アノテータ処理部104は、「××県××市××〇

10

20

30

40

50

丁目○番地○号に」を新たな候補文字列としてもよい。

【0044】

一方、ステップS405において、ジオコーディングによる位置情報の取得に失敗したと判断された場合（ステップS405：No）、位置情報の取得に成功した候補文字列がメモリに格納されているか否かを判断する（ステップS409）。

【0045】

ここで、「ジオコーディングによる位置情報の取得に失敗した場合」とは、実際に、ジオコーディングのAPIにてエラーが発生した場合に限らず、例えば、候補文字列の一部の文字列から位置情報が取得された場合（すなわち、候補文字列が完全一致しなかった場合）も含む。例えば、APIによっては、「××県××市××○丁目○番地○号に」を入力した場合に、エラーが発生せずに、位置情報の取得が可能な一部の文字列「××県××市××○丁目○番地○号」から、位置情報を取得する場合がある。この場合、アノテータ処理部104は、「ジオコーディングによる位置情報の取得に失敗した」と判断するようにしてもよい。

10

【0046】

また、APIによっては、候補文字列との一致度を示す信頼度を返す場合がある。この場合、例えば、アノテータ処理部104は、直前の候補文字列から信頼度が上昇した場合または直前の候補文字列と信頼度が同一の場合、「ジオコーディングによる位置情報の取得に成功した」と判断し、直前の候補文字列から信頼度が低下した場合、「ジオコーディングによる位置情報の取得に失敗した」と判断するようにしてもよい。

20

【0047】

ステップS409において、位置情報の取得に成功した候補文字列がメモリに格納されていないと判断された場合（ステップS409：No）、アノテータ処理部104は、図4に示す一連の処理を終了する。

【0048】

一方、ステップS409において、位置情報の取得に成功した候補文字列がメモリに格納されていると判断された場合（ステップS409：Yes）、アノテータ処理部104は、メモリに格納されている候補文字列が、区、地番、および枝番を含むか否かを判断する（ステップS410）。ステップS410において、メモリに格納されている候補文字列が、区、地番、および枝番を含まないと判断された場合（ステップS410：No）、アノテータ処理部104は、図4に示す一連の処理を終了する。

30

【0049】

一方、ステップS410において、メモリに格納されている候補文字列が、区、地番、および枝番を含むと判断された場合（ステップS410：Yes）、アノテータ処理部104は、メモリに格納されている候補文字列の長さが、所定文字数未満であるか否かを判断する（ステップS411）。ステップS411において、メモリに格納されている候補文字列の長さが、所定文字数未満ではないと判断された場合（ステップS411：No）、アノテータ処理部104は、図4に示す一連の処理を終了する。

【0050】

一方、ステップS411において、メモリに格納されている候補文字列の長さが、所定文字数未満であると判断された場合（ステップS411：Yes）、アノテータ処理部104は、メモリに格納されている候補文字列を、住所とする文字列に決定する（ステップS412）。そして、アノテータ処理部104は、図4に示す一連の処理を終了する。

40

【0051】

（アノテータ処理部104によるナイーブベイズ推定値取得処理の手順）

図5は、実施形態に係るアノテータ処理部104によるナイーブベイズ推定値取得処理の手順を示すフローチャートである。図5は、図3にフローチャートにおけるステップS307の処理を詳細に説明するものである。なお、図5の処理には、複数のPOI名称の候補文字列を含む候補リストが入力される。これに応じて、アノテータ処理部104は、図5の処理を、候補リストに含まれるPOI名称の候補文字列毎に実行する。

50

## 【 0 0 5 2 】

まず、アノテータ処理部 1 0 4 は、P O I 名称の候補文字列をナイーブベイズ分類器にかけることにより、P O I 名称としての確からしさの推定を行う（ステップ S 5 0 1）。そして、アノテータ処理部 1 0 4 は、ステップ S 5 0 1 で推定された P O I 名称としての確からしさを示す、ナイーブベイズ推定値を取得する（ステップ S 5 0 2）。このナイーブベイズ推定値は、「 1 . 0 0 」を最大値とするものである。

## 【 0 0 5 3 】

次に、アノテータ処理部 1 0 4 は、P O I 名称の候補文字列が、所定文字列を含むか、または、強調されているか否かを判断する（ステップ S 5 0 3）。所定文字列とは、例えば、「株式会社」、「店」等、P O I 名称である可能性を高める文字列である。このため、情報抽出装置 1 0 0 には、予め、複数の所定の文字列が、メモリ等に予め登録されている。また、P O I 名称の候補文字列が強調されている場合とは、例えば、P O I 名称の候補文字列に強調タグ（例えば、<strong>、<em>、<b>、<font>等）が付されている場合等である。

## 【 0 0 5 4 】

ステップ S 5 0 3 において、P O I 名称の候補文字列が、所定文字列を含まず、且つ、強調されていないと判断された場合（ステップ S 5 0 3 : N o）、アノテータ処理部 1 0 4 は、ステップ S 5 0 5 へ処理を進める。

## 【 0 0 5 5 】

一方、ステップ S 5 0 3 において、P O I 名称の候補文字列が、所定文字列を含む、または、強調されていると判断された場合（ステップ S 5 0 3 : Y e s）、アノテータ処理部 1 0 4 は、ステップ S 5 0 2 で取得されたナイーブベイズ推定値を所定倍（例えば、1 . 2 5 倍）する（ステップ S 5 0 4）。但し、これによりナイーブベイズ推定値が「 1 . 0 0 」を超える場合は、ナイーブベイズ推定値を「 1 . 0 0 」とする。そして、アノテータ処理部 1 0 4 は、ステップ S 5 0 5 へ処理を進める。

## 【 0 0 5 6 】

ステップ S 5 0 5 では、アノテータ処理部 1 0 4 は、ナイーブベイズ推定値を、P O I 名称の候補文字列に対応付けて、候補リストに保存する。そして、アノテータ処理部 1 0 4 は、図 5 に示す一連の処理を終了する。

## 【 0 0 5 7 】

以上説明したように、本実施形態に係る情報抽出装置 1 0 0 によれば、W E B ページ 1 1 0 から抽出した候補文字列に対してジオコーディングを行うことにより、位置情報の取得を試み、当該位置情報の取得結果に応じて、住所とする文字列を決定することができる。特に、本実施形態に係る情報抽出装置 1 0 0 によれば、位置情報の取得に成功した場合、位置情報の再取得に失敗するまで、候補文字列を後方に延長して、位置情報の再取得を試みることができる。これにより、本実施形態に係る情報抽出装置 1 0 0 によれば、W E B ページ 1 1 0 から抽出する住所文字列の終端を高精度に特定することができる。したがって、本実施形態に係る情報抽出装置 1 0 0 によれば、W E B ページ 1 1 0 からの P O I データの抽出精度を高めることができる。

## 【 0 0 5 8 】

また、本実施形態に係る情報抽出装置 1 0 0 によれば、W E B ページ 1 1 0 から抽出した P O I 名称の候補文字列が所定文字列を含むか、または、強調されている場合、その P O I 名称としての確からしさを示すナイーブベイズ推定値を高めることができる。このため、本実施形態に係る情報抽出装置 1 0 0 によれば、例えば、W E B ページ 1 1 0 から複数の P O I 名称の候補文字列が抽出された場合であっても、これら複数の P O I 名称の候補文字列の中から、ナイーブベイズ推定値に基づいて、実際の P O I 名称である可能性が最も高い文字列を抽出することができる。したがって、本実施形態に係る情報抽出装置 1 0 0 によれば、W E B ページ 1 1 0 からの P O I データの抽出精度を高めることができる。

## 【 0 0 5 9 】



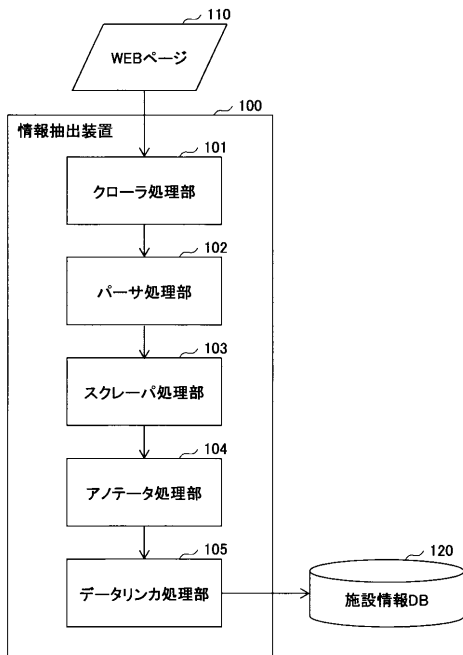
以上、本発明の好ましい実施形態について詳述したが、本発明はこれらの実施形態に限定されるものではなく、特許請求の範囲に記載された本発明の要旨の範囲内において、種々の変形又は変更が可能である。

【符号の説明】

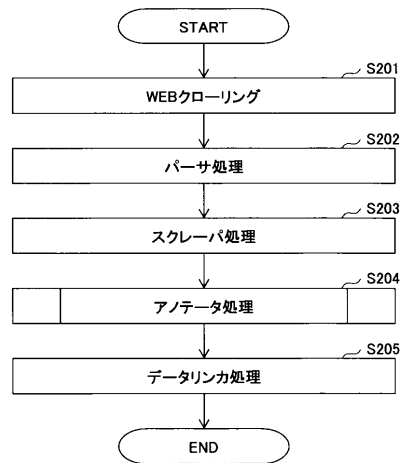
【0060】

- 100 情報抽出装置
- 101 クローラ処理部（文書データ取得部）
- 102 パーサ処理部
- 103 スクレーパー処理部
- 104 アノテータ処理部（候補文字列抽出部、位置情報取得部、住所文字列決定部）
- 105 データリンカ処理部
- 110 WEBページ
- 120 施設情報DB

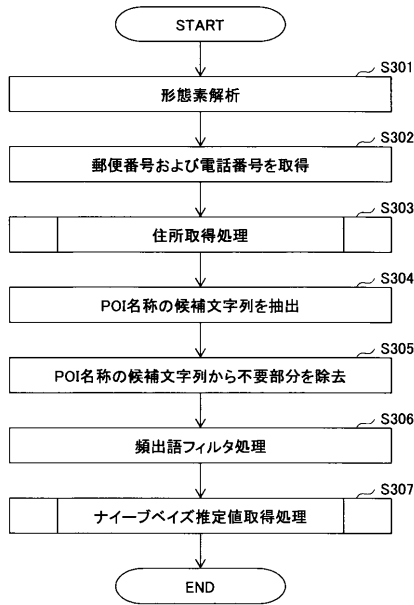
【図1】



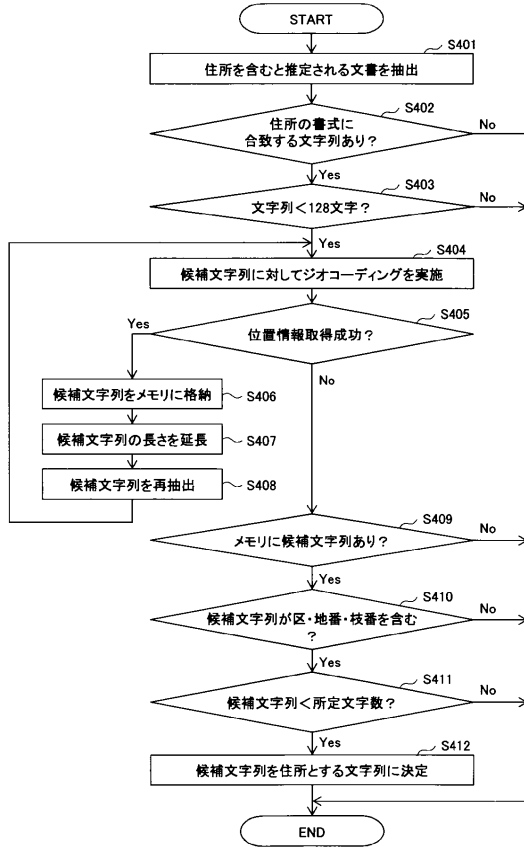
【図2】



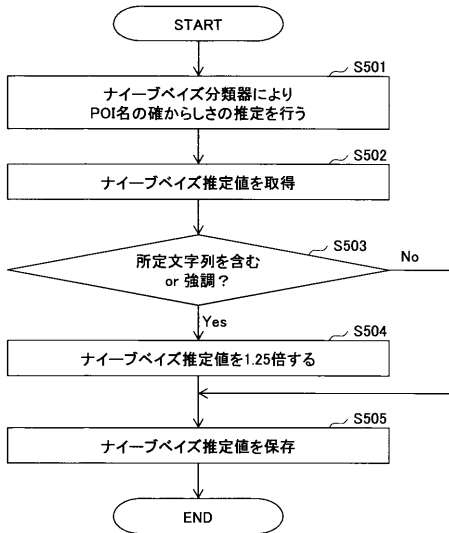
【図3】



【図4】



【図5】



## フロントページの続き

- (56)参考文献 特表2009-506459(JP,A)  
特開2004-086272(JP,A)  
特開2007-179329(JP,A)  
特開2016-091315(JP,A)  
特開平10-240710(JP,A)  
特開2006-064443(JP,A)  
福田拓真, 力宗幸男, Geoマイクロフォーマットを用いた住所自動検出・地図表示システムの開発, 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 2010年 1月14日, 第109巻, 第379号, p. 93~98  
石田武久, 外2名, Web上のイラストマップを実地図に重ね合わせるシステムの試作, 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 2011年12月 9日, 第111巻, 第361号, p. 43~48

## (58)調査した分野(Int.Cl., DB名)

G06F 16/00 - 16/958  
G06F 40/00 - 40/197