

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局



(43) 国際公開日
2011年9月1日(01.09.2011)

PCT

(10) 国際公開番号
WO 2011/105606 A1

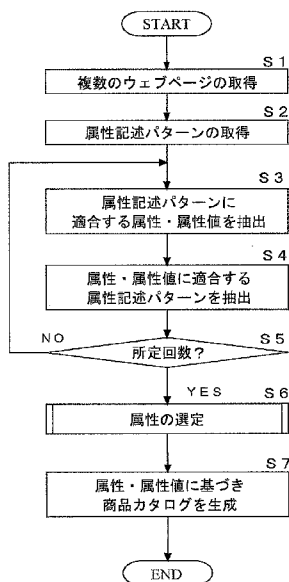
- (51) 国際特許分類:
G06F 17/30 (2006.01) G06Q 30/00 (2006.01)
- (21) 国際出願番号: PCT/JP2011/054510
- (22) 国際出願日: 2011年2月28日(28.02.2011)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:
特願 2010-043392 2010年2月26日(26.02.2010) JP
特願 2010-043391 2010年2月26日(26.02.2010) JP
特願 2010-043390 2010年2月26日(26.02.2010) JP
- (71) 出願人(米国を除く全ての指定国について): 楽天株式会社(Rakuten, Inc.) [JP/JP]; 〒1400002 東京都品川区東品川四丁目12番3号 Tokyo (JP).
- (72) 発明者; および
- (75) 発明者/出願人(米国についてのみ): 関根 聡 (SEKINE Satoshi) [JP/JP]; 〒1520031 東京都目黒区中根一丁目6番22号 株式会社ランゲージ・クラフト研究所内 Tokyo (JP). 竹中 孝真 (TAKENAKA Takamasa) [JP/JP]; 〒1400002 東京都品川区東品川四丁目12番3号 楽天株式会社内 Tokyo (JP).
- (74) 代理人: 特許業務法人 インテクト国際特許事務所, 外 (INTECT INTERNATIONAL PATENT OFFICE et al.); 〒1020083 東京都千代田区麹町四丁目7番2号 サンライン第7ビル4階 Tokyo (JP).
- (81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PE, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.
- (84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), ユーラシア

[続葉有]

(54) Title: INFORMATION PROCESSING DEVICE, INFORMATION PROCESSING METHOD, PROGRAM FOR INFORMATION PROCESSING DEVICE, AND RECORDING MEDIUM

(54) 発明の名称: 情報処理装置、情報処理方法、情報処理装置用のプログラム、および、記録媒体

[図4]



S1 Acquisition of plurality of web pages
 S2 Acquisition of attribute description patterns
 S3 Extraction of attributes/attribute values compatible with attribute description patterns
 S4 Extraction of attribute description patterns compatible with attributes/attribute values
 S5 Predetermined number?
 S6 Selection of attributes
 S7 Generation of product catalogue on basis of attributes/attribute values

(57) Abstract: The disclosed information processing device: obtains a plurality of web pages from the same category, said category classifying an object listed in a web page (S1); obtains attribute-related terms, which are related to the attributes of the object listed in the web pages, or attribute description patterns, which are used to describe the attributes of the object, as initial data (S2); extracts attribute-related terms for attributes which are compatible with the attribute description patterns from the plurality of web pages (S3); and extracts attribute description patterns which are compatible with the attribute-related words from the plurality of web pages (S4).

(57) 要約: ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得し (S1)、初期データとして、ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得し (S2)、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出し (S3)、複数のウェブページから、属性関連語に適合する属性記述パターンを抽出する (S4)。



WO 2011/105606 A1

(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), ヨーロッパ 添付公開書類:

(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI
— 国際調査報告 (条約第 21 条(3))
(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

明 細 書

発明の名称：

情報処理装置、情報処理方法、情報処理装置用のプログラム、および、記録媒体

技術分野

[0001] 本発明は、インターネット上のウェブページを分析する情報処理装置、情報処理方法、情報処理装置用のプログラム、および、記録媒体の技術分野に関する。

背景技術

[0002] インターネット上のウェブサイトには商品等を扱ったウェブページが多数あり、ユーザはそのウェブページを閲覧して商品の購入を行っている。閲覧して購入する際、ユーザは通常、検索を行っている。この検索の技術において、多くのウェブページの中からユーザが欲しい商品を素早く探すために検索の技術の向上を図ることが行われている。例えば、特許文献1には、商品のカテゴリ別検索を順に大分類、中分類、小分類と検索する3層のカテゴリ別検索に限定して、6桁の整数分類コード表作成方法を考案し、この分類コード表作成方法を使用して商品分類コード表及び店舗分類コード表を作成し、これらの分類コード表をインターネットショッピングモールに設け、ショップが簡単に商品及び店舗情報の登録ができユーザが簡単に商品及び店舗を検索ができるショッピングモールにおける検索システムが開示されている。

先行技術文献

特許文献

[0003] 特許文献1：特開2002-236694号公報

発明の概要

発明が解決しようとする課題

[0004] ところで、ワインのような商品の場合、産地、容量といった商品の属性に

関する情報や、旅行関連サービスのようなサービスの場合、料金、アクセスといったサービスの属性に関する情報が、ウェブページの中に記載されている。このような商品等の属性を抽出すれば、さまざまなサービスを提供できる可能性がある。しかし、特許文献1のような従来では、商品等の属性を抽出することは難しかった。

[0005] 本発明は、このような問題に鑑みてなされたものであり、その課題の一例は、ウェブページから商品等の属性を自動的に取得する情報処理装置等を提供することを目的とする。

課題を解決するための手段

[0006] 上記課題を解決するために、請求項1に記載の発明は、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段と、初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段と、前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段と、前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段と、を備えたことを特徴とする。

[0007] 請求項2に記載の発明は、請求項1に記載の情報処理装置において、前記属性抽出手段および前記属性記述パターン抽出手段を交互に繰り返す繰返手段を更に備えたことを特徴とする。

[0008] 請求項3に記載の発明は、請求項1または請求項2に記載の情報処理装置において、前記属性抽出手段が、前記属性関連語として、前記属性の属性名を抽出することを特徴とする。

[0009] 請求項4に記載の発明は、請求項1から請求項3のいずれか1項に記載の情報処理装置において、抽出された前記属性関連語から属性リストを生成する属性リスト生成手段と、抽出された前記属性記述パターンのパターンリストを生成するパターンリスト生成手段と、を更に備えたことを特徴とする。

- [0010] 請求項 5 に記載の発明は、請求項 1 から請求項 4 のいずれか 1 項に記載の情報処理装置において、前記属性関連語のスコア付けを行う属性スコアリング手段と、前記スコアの順に前記属性関連語のランク付けを行い、所定のランク以上の属性関連語を選択する属性選択手段と、を更に備えたことを特徴とする。
- [0011] 請求項 6 に記載の発明は、請求項 5 に記載の情報処理装置において、前記属性スコアリング手段が、前記属性関連語の検索のヒット件数に基づき、前記属性関連語のスコア付けを行うことを特徴とする。
- [0012] 請求項 7 に記載の発明は、請求項 5 に記載の情報処理装置において、前記属性スコアリング手段が、前記対象を販売する複数の店舗を有するウェブサイトにおいて、前記属性関連語が出現しているウェブページの前記店舗の数に基づき、前記属性関連語のスコア付けを行うことを特徴とする。
- [0013] 請求項 8 に記載の発明は、請求項 1 から請求項 7 のいずれか 1 項に記載の情報処理装置において、前記カテゴリとは異なるカテゴリに属している複数のウェブページにおいて出現する前記属性関連語を取り除く属性フィルタ手段を更に備えたことを特徴とする。
- [0014] 請求項 9 に記載の発明は、請求項 1 から請求項 8 のいずれか 1 項に記載の情報処理装置において、前記属性記述パターンのスコア付けを行う属性記述パターン・スコアリング手段と、前記スコアの順に前記属性記述パターンのランク付けを行い、所定のランク以上の属性記述パターンを選択する属性記述パターン選択手段と、を更に備えたことを特徴とする。
- [0015] 請求項 10 に記載の発明は、請求項 9 に記載の情報処理装置において、前記属性記述パターン・スコアリング手段が、前記属性関連語と前記属性記述パターンとが共に出現する共起数に基づき前記属性記述パターンのスコア付けを行うことを特徴とする。
- [0016] 請求項 11 に記載の発明は、請求項 1 から請求項 10 のいずれか 1 項に記載の情報処理装置において、前記属性名同士が類似であるか否かを判定する属性名類似判定手段と、前記属性名類似判定手段により類似と判定された属

性を集約する属性名集約手段と、を更に備えたことを特徴とする。

- [0017] 請求項 1 2 に記載の発明は、請求項 1 1 に記載の情報処理装置において、前記属性抽出手段が、前記属性関連語として、前記属性名および前記属性名に対応する属性値を抽出し、前記属性名集約手段が、前記属性値に基づき前記属性名を集約することを特徴とする。
- [0018] 請求項 1 3 に記載の発明は、請求項 1 から請求項 1 0 のいずれか 1 項に記載の情報処理装置において、前記ウェブページ取得手段が、前記対象の供給元のウェブページを取得し、前記初期データ取得手段、前記属性抽出手段、および、前記属性記述パターン抽出手段により、前記対象の供給元のウェブページから供給元対象属性関連語を抽出し、前記供給元対象属性関連語と前記属性関連語とを比較する属性関連語比較手段を更に備えたことを特徴とする。
- [0019] 請求項 1 4 に記載の発明は、請求項 1 から請求項 1 3 のいずれか 1 項に記載の情報処理装置において、抽出された前記属性関連語に基づき、前記属性関連語が記載されたカタログを生成するカタログ生成手段を更に備えたことを特徴とする。
- [0020] 請求項 1 5 に記載の発明は、請求項 1 から請求項 1 4 のいずれか 1 項に記載の情報処理装置において、前記複数のウェブページから、前記属性関連語の出現回数が所定回数以下のウェブページを抽出するウェブページ抽出手段と、
を更に備えたことを特徴とする。
- [0021] 請求項 1 6 に記載の発明は、請求項 1 5 に記載の情報処理装置において、前記ウェブページ抽出手段が、前記属性関連語の出現回数がゼロのウェブページを抽出することを特徴とする。
- [0022] 請求項 1 7 に記載の発明は、請求項 1 から請求項 1 6 のいずれか 1 項に記載の情報処理装置において、前記属性関連語に基づき、前記複数のウェブページをグルーピングするウェブページ・属性グルーピング手段を更に備えたことを特徴とする。

- [0023] 請求項 18 に記載の発明は、情報処理装置が情報処理をする情報処理方法において、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得ステップと、前記ウェブページに記載されている対象の属性の記述に用いられる属性記述パターンを取得する属性記述パターン取得ステップと、前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出ステップと、抽出された前記属性関連語に基づき、前記属性抽出ステップで使用する前記属性記述パターンを、前記複数のウェブページから、更に抽出する属性記述パターン抽出ステップと、を有することを特徴とする。
- [0024] 請求項 19 に記載の発明は、情報処理装置が情報処理をする情報処理方法において、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得ステップと、前記ウェブページに記載されている対象の属性に関連した属性関連語を取得する属性関連語取得ステップと、前記属性の記述に用いられる属性記述パターンであって、前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出ステップと、抽出された前記属性関連語に基づき、前記属性記述パターン抽出ステップで使用する属性関連語を、前記複数のウェブページから、更に抽出する属性関連語抽出ステップと、を有することを特徴とする。
- [0025] 請求項 20 に記載の発明は、コンピュータを、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段、初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段、前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段、および、前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段として機能させることを特徴とする。

[0026] 請求項 2 1 に記載の発明は、コンピュータを、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段、初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段、前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段、および、前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段として機能させることを特徴とする情報処理装置用のプログラムを記録する。

発明の効果

[0027] 本発明によれば、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得し、初期データとして、ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得し、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出し、複数のウェブページから、属性関連語に適合する属性記述パターンを抽出することにより、同一のカテゴリに属している複数のウェブページから、属性関連語を抽出し、属性記述パターンを抽出するか、または、属性記述パターンを抽出し、属性関連語を抽出しているため、同一のカテゴリに含まれる属性を精度良く取得できる。

図面の簡単な説明

[0028] [図1]本発明の第 1 実施形態に係る情報処理システムの概要構成例を示す模式図である。

[図2]図 1 の情報処理サーバの概要構成の一例を示すブロック図である。

[図3]図 1 のショッピングサーバの概要構成の一例を示すブロック図である。

[図4]図 1 の情報処理サーバにおいてカタログを生成する動作例を示すフローチャートである。

[図5]図 1 のショッピングサーバのウェブページの一例を示す説明図である。

[図6] 図5のウェブページのソースコードの一例を示す説明図である。

[図7] 図2の属性記述パターンデータベースに記憶された属性記述パターンの一例を示す模式図である。

[図8] 属性・属性値の抽出の様子の一列を示す模式図である。

[図9] 抽出された属性・属性値の一例を示す模式図である。

[図10] 生成された商品等カタログの一例を示す模式図である。

[図11] 図1の情報処理サーバにおける属性選定のサブルーチンの一例を示すフローチャートである。

[図12] 図4の商品等のカタログ生成の第1変形例の動作例を示すフローチャートである。

[図13] 図12の第1変形例の属性・属性値の抽出の様子の一列を示す模式図である。

[図14] 図4の商品等のカタログ生成の第2変形例の動作例を示すフローチャートである。

[図15] 図14の属性・属性値抽出のサブルーチンの一例を示すフローチャートである。

[図16] 図14の属性記述パターン抽出のサブルーチンの一例を示すフローチャートである。

[図17] 図1の情報処理サーバにおいて属性・属性値を判定する動作例を示すフローチャートである。

[図18] 商品等供給元のウェブページの一例を示す説明図である。

[図19] 生成された商品等カタログの一例を示す模式図である。

[図20] 第2実施形態に係る情報処理システムにおいてウェブページを抽出する動作例を示すフローチャートである。

[図21] 図20のウェブページ抽出の第1変形例の動作例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

[0029] 以下、図面を参照して本発明の実施形態について説明する。

(第1実施形態)

[0030] [1. 情報処理システムの構成および機能概要]

まず、本発明の第1実施形態に係る情報処理システムの構成および概要機能について、図1を用いて説明する。

[0031] 図1は、本実施形態に係る情報処理システム1の概要構成例を示す模式図である。

[0032] 図1に示すように、情報処理システム1は、ウェブページから商品等のカタログ生成したり、誤ったカテゴリに登録されたウェブページを抽出したりする情報処理サーバ（情報処理装置の一例）10と、ショッピングサイトを運営するためや、ブログのため情報提供サーバ20と、ショッピングサイトに出店している店舗主の店舗主端末30と、ショッピングサイトで商品等（ショッピングサイトで提供されているサービスを含む）を購入したり、ブログを投稿するユーザのユーザ端末35と、を備えている。なお、商品等や、ブログは、ウェブページに記載されている対象の一例である。

[0033] 情報処理サーバ10と、情報提供サーバ20とは、ローカルエリアネットワーク等により接続され、相互にデータの送受信が可能になっていて、サーバシステム5を構成している。そして、サーバシステム5と、店舗主端末30と、ユーザ端末35とは、ネットワーク3により接続され、例えば、通信プロトコルにTCP/IP等を用いて相互にデータの送受信が可能になっている。なお、ネットワーク3は、例えば、インターネット、専用通信回線（例えば、CATV（Community Antenna Television）回線）、移動体通信網（基地局等を含む）、およびゲートウェイ等により構築されている。

[0034] 情報処理システム1は、ウェブページからカタログを生成するカタログ生成システムとして、または、誤ったカテゴリに登録されたウェブページを抽出するウェブページ抽出システムとして機能する。

[0035] 情報処理サーバ10は、情報提供サーバ20等に登録されたウェブページから商品等のカタログを生成したり、当該カタログをユーザ端末35等から閲覧できるようにしたりする。また、情報処理サーバ10は、情報提供サー

バ20等に登録されたウェブページから誤ったカテゴリに登録されたウェブページを抽出したり、抽出結果に基づき、情報提供サーバ20上のウェブページの整理を行ったり、ウェブページを登録した店舗主等の店舗主端末30等に抽出結果を通知したりする。

[0036] 情報提供サーバ20は、ショッピングサイトで商品等を販売するためのウェブサーバや、データベースサーバ等として機能し、ウェブページの登録の受け付けや、ユーザ登録や、商品等の購入手続き等の各種処理を行う。また、情報提供サーバ20は、商品等のカテゴリ毎に分類されたウェブページを有している。また、情報提供サーバ20は、ユーザからのブログの投稿を受け付け、ブログの内容等に基づきカテゴリ毎に分類しインターネット上に公開する。

[0037] 店舗主が使用する店舗主端末30は、パーソナルコンピュータや携帯型無線電話機やPDA (Personal Digital Assistant) 等の携帯端末である。店舗主は店舗主端末30を使用して、ウェブページを情報提供サーバ20に登録したり、更新したりする。

[0038] ユーザが使用するユーザ端末35は、パーソナルコンピュータや携帯型無線電話機やPDA等の携帯端末である。ユーザはユーザ端末35を使用して、商品等の検索や商品等の購入等を行う。

[0039] [2. 各サーバの構成および機能]

(2. 1 情報処理サーバ10の構成および機能)

次に、情報処理サーバ10の構成および機能について、図2を用いて説明する。

[0040] 図2は、情報処理サーバ10の概要構成の一例を示すブロック図である。

[0041] 図2に示すように、コンピュータとして機能する情報処理サーバ10は、通信部11と、記憶部12と、入出力インターフェース部13と、システム制御部14と、を備えている。そして、システム制御部14と入出力インターフェース部13とは、システムバス15を介して接続されている。

[0042] 通信部11は、ネットワーク3に接続してユーザ端末35等との通信状態

を制御したり、ローカルエリアネットワークに接続して、情報提供サーバ20等の他のサーバとデータの送受信を行ったりする。

[0043] 記憶部12は、例えば、ハードディスクドライブ等により構成されており、オペレーティングシステムおよびサーバプログラム等の各種プログラムや、データ等を記憶する。なお、各種プログラムは、例えば、他のサーバ装置等からネットワーク3を介して取得されるようにしてもよいし、記録媒体に記録されてドライブ装置を介して読み込まれるようにしてもよい。

[0044] また、記憶部12には、属性記述パターンデータベース（以下「属性記述パターンDB」とする。）12a、属性・属性値データベース（以下「属性・属性値DB」とする。）12b等が構築されている。

[0045] 属性記述パターンDB12aには、商品等やブログの属性の記述に用いられる属性記述パターンの初期データや、ウェブページから抽出した属性記述パターンが記憶されている。なおブログの属性としてブログのカテゴリが挙げられる。

[0046] 属性・属性値DB12bには、情報処理サーバ10による処理の一例として、ウェブページから抽出した商品等の属性に関する属性名と属性値とが記憶される。ここで、属性関連語の一例として、属性名のみや、属性名を含む語句や、属性名と属性値との組等が挙げられる。また、属性・属性値という表記は、属性と属性値とが対になっている場合で、具体的に属性名と属性値とが組になった場合も含む。

[0047] 次に、入出インターフェース部13は、通信部11および記憶部12とシステム制御部14との間のインターフェース処理を行うようになっている。

[0048] システム制御部14は、CPU（Central Processing Unit）14a、ROM（Read Only Memory）14b、RAM（Random Access Memory）14c等により構成されている。そして、システム制御部14は、CPU14aが、ROM14bや記憶部12に記憶された各種プログラムを読み出し実行することにより、複数のウェブページから、属性記述パターンに適合する属性名

や属性値を抽出したり、抽出した属性名や属性値から商品等のカタログを生成したりする。また、システム制御部 14 は、複数のウェブページから、属性の属性名の出現回数が所定回数以下のウェブページを、誤ったカテゴリに登録されたウェブページとして抽出したりする。

[0049] (2. 2 情報提供サーバ 20 の構成および機能)

次に、情報提供サーバ 20 の構成および機能について、図 3 を用いて説明する。

図 3 は、情報提供サーバ 20 の概要構成の一例を示すブロック図である。

[0050] 図 3 に示すように、情報提供サーバ 20 は、通信部 21 と、記憶部 22 と、入出インターフェース部 23 と、システム制御部 24 と、を備え、システム制御部 24 と入出インターフェース部 23 とは、システムバス 25 を介して接続されている。なお、情報提供サーバ 20 の構成および機能は、情報処理サーバ 10 の構成および機能とほぼ同じであるので、情報処理サーバ 10 の各構成や各機能において、異なるところを中心に説明する。

[0051] 通信部 21 は、ネットワーク 3 やローカルエリアネットワーク等を通して、店舗主端末 30 やユーザ端末 35 や情報処理サーバ 10 等と通信状態を制御等するようになっている。

[0052] 記憶部 22 には、商品データベース（以下「情報 DB」とする。） 22 a や、会員データベース（以下「会員 DB」とする。） 22 b や商品等カタログデータベース（以下「商品等カタログ DB」とする。） 22 c 等が構築されている。

[0053] 情報 DB 22 a には、ウェブページに記載されている対象の一例である商品、サービス、および、ブログ等に関する情報が記憶されている。例えば、情報 DB 22 a には、商品等を識別するための識別子である商品 ID に関連付けられ、商品名（サービス名を含む）、種類、商品の画像、サービスに関連した画像、スペック、および、商品等の紹介の要約文等の商品情報や、広告情報等が記憶されている。また、情報 DB 22 a には、ユーザが投稿してきたブログの記事がカテゴリ分けされて記憶されている。また、情報 DB 2

2 aには、HTML (HyperText Markup Language)、XML (Extensible Markup Language) 等のマークアップ言語等により記述されたウェブページのファイル等が記憶されている。また、情報DB 2 2 aには、製造元情報（製造元ドメインを含む）および販売元情報（販売元ドメインを含む）等の商品供給元の情報が記憶されていて、各商品の商品IDに、各商品の公式の情報が記載されている商品供給元のウェブページのURL (Uniform Resource Locator) 等が関連づけられている。

[0054] 会員DB 2 2 bには、会員登録されたユーザ（インターネットショップの利用者）のユーザID、名称、住所、電話番号、メールアドレス等のユーザ情報が登録されている。このようなユーザ情報は、ユーザIDによってユーザ毎に判別可能になっている。また、会員DB 2 2 bには、ユーザがユーザ端末3 5からインターネットショップのサイトにログインする際に必要な、ユーザID、ログインID、および、パスワードが登録されている。ここで、ログインIDおよびパスワードは、ログイン処理（ユーザの認証処理）に使用されるログイン情報である。

[0055] 商品等カタログDB 2 2 cには、情報処理サーバ1 0により生成された商品等カタログが商品カテゴリ毎、商品毎に記憶される。

[0056] システム制御部2 4は、CPU 2 4 a、ROM 2 4 b、RAM 2 4 c等により構成されている。そして、システム制御部2 4は、CPU 2 4 aが、ROM 2 4 bや記憶部2 2に記憶された各種プログラムを読み出し実行することにより、店舗主によるウェブページの登録や更新や、ユーザによる商品購入処理や、商品の購買履歴をユーザID毎に記録させたりする。またユーザ端末3 5からの要求により、商品等カタログの情報を送信したりする。

[0057] [3. 第1実施形態の商品等のカタログ生成システムの動作]

次に、本発明の一実施形態に係る情報処理システム1のカタログ生成システムとしての動作について図4～図11を用いて説明する。

[0058] 図4は、情報処理サーバ1 0においてウェブページを抽出する動作例を示すフローチャートである。図5は、情報提供サーバ2 0のウェブページの一

例を示す説明図である。図6は、ウェブページのソースコードの一例を示す説明図である。図7は、属性記述パターンDBに記憶された属性記述パターンの一例を示す模式図である。図8は、属性・属性値の抽出の様子の一列を示す模式図である。図9は、抽出された属性・属性値の一例を示す模式図である。図10は、生成された商品等カタログの一例を示す模式図である。図11は、情報処理サーバ10における属性選定のサブルーチンの一例を示すフローチャートである。

[0059] (3. 1. 商品等のカタログ生成の流れ)

まず、商品等のカタログ生成の流れについて、図4を用いて説明する。

[0060] 図4に示すように、情報処理サーバ10は、複数のウェブページを取得する(ステップS1)。具体的には、情報処理サーバ10のシステム制御部14は、情報提供サーバ20により運営されるショッピングサイトの同一のカテゴリに所属している商品に関して、このカテゴリの全ウェブページを、通信部11を通して情報DB22aから取得する。さらに具体的には、図5に示すように、テキスト部51、52、53、54のテキストデータを含むウェブページ50等が取得される。また、ウェブページ50のソースコードは、図6に示すように、HTML等のマークアップ言語等で記述されている。このように、情報処理サーバ10のシステム制御部14および通信部11は、ウェブページに記載されている対象を分類するカテゴリにおいて、同一のカテゴリに属している複数のウェブページを取得するウェブページ取得手段の一例として機能する。

[0061] 次に、情報処理サーバ10は、属性記述パターンを取得する(ステップS2)。具体的には、情報処理サーバ10のシステム制御部14は、下記のステップS3~S5におけるブートストラップ法の初期データとして、図7に示すように、属性記述パターンDB12aの属性記述パターンリストから、初期の属性記述パターンを取得する。ここで、属性記述パターンは、図7に示すように、前部、中部、および、後部に分かれていて、例えば、属性記述パターン” [:] ” の場合、前部” [”、中部” : ”、および、後部”

] ”である。前部と中部との間の語句が属性名で、中部と後部との間の語句が属性値である。また、属性記述パターンには、HTMLタグの要素が含まれる場合がある。このように情報処理サーバ10のシステム制御部14は、初期データとして、ウェブページに記載されている対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段の一例として機能する。

[0062] 次に、情報処理サーバ10は、属性記述パターンに適合する属性・属性値を抽出する（ステップS3）。具体的には、情報処理サーバ10のシステム制御部14は、ウェブページ50等の複数のウェブページの中から、図8に示すように、属性記述パターン61等に適合する語句の部分（例えば” [品種：○○○] ”）を取り出し、属性名” 品種” や、属性名” 品種” に対応した属性値” ○○○” 等を抽出する。そして、抽出した属性名および属性値は、属性リストとして属性・属性値DB12bに記憶される。ここで、どんなパターンにもマッチする特殊文字、すなわち、” * ” や ” ? ” 等のワイルドカードと属性記述パターンとが用いられて、属性・属性値が抽出される。なお、属性・属性値の例として、旅行関連サービスの場合、[宿泊料金：○○○]、ブログであるイベント紹介をしている場合、[会場：○○○] 等が挙げられる。

[0063] このように情報処理サーバ10のシステム制御部14は、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出する属性抽出手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、属性関連語として、属性の属性名を抽出する属性抽出手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、抽出された属性関連語から属性リストを生成する属性リスト生成手段として機能する。

[0064] 次に、情報処理サーバ10は、属性・属性値に適合する属性記述パターンを抽出する（ステップS4）。具体的には、情報処理サーバ10のシステム制御部14は、図8に示すように、属性・属性値62（例えば、属性名” 品

種” および属性値” ○○○”) に適合する (例えば、<td> 品種</td><td>○○○○</td>) を取り出し、属性記述パターンをウェブページ 50 等の複数のウェブページの中から抽出する。そして、抽出した属性記述パターンは、図 7 に示すように、属性記述パターンリストに追加され、属性記述パターン DB 12a に記憶される。なお、例えば、” 容量 *ml” のように、属性値に関しては、抽出された属性値自体でなく、ワイルドカードが用いられてもよい。

[0065] このように情報処理サーバ 10 のシステム制御部 14 は、複数のウェブページから、属性関連語に適合する属性記述パターンを抽出する属性記述パターン抽出手段の一例として機能する。また、情報処理サーバ 10 のシステム制御部 14 は、抽出された属性記述パターンのパターンリストを生成するパターンリスト生成手段として機能する。

[0066] 次に、情報処理サーバ 10 は、所定回数を判定する (ステップ S5)。具体的には、情報処理サーバ 10 のシステム制御部 14 は、ステップ S3 およびステップ S4 を反復実行した回数が所定回数に達しているか否かを判定する。そして、所定回数に達していない場合 (ステップ S5 ; NO)、情報処理サーバ 10 のシステム制御部 14 は、ステップ S3 に戻り、抽出した新たな属性記述パターンにより、新たな属性・属性値を抽出する。情報処理サーバ 10 のシステム制御部 14 は、所定回数に達するまで、ステップ S3 およびステップ S4 を繰り返す。

[0067] このように情報処理サーバ 10 のシステム制御部 14 は、ステップ S2 からステップ S4 において、商品の属性の記述に用いられる属性記述パターンを取得する属性記述パターン取得ステップと、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出する属性抽出ステップと、抽出された属性関連語に基づき、属性抽出ステップで使用する属性記述パターンを、複数のウェブページから更に抽出する属性記述パターン抽出ステップとを実行する。情報処理サーバ 10 のシステム制御部 14 は、属性抽出手段および属性記述パターン抽出手段を交互に繰り返す繰返手段の一例とし

て機能する。

[0068] 所定回数に達した場合（ステップS 5；YES）、情報処理サーバ10は、属性の選定を行う（ステップS 6）。具体的には、情報処理サーバ10のシステム制御部14は、ステップS 3で抽出した属性名および属性値から、属性選定のサブルーチンにより属性を選定する。属性選定のサブルーチンでは、情報処理サーバ10のシステム制御部14は、属性にスコアを付けてランク付けしたり、ノイズの属性を除去したり、同義語の属性を集約する（詳細は後述）。図9に示すように、カテゴリ”ワイン”においては、属性名”品種”、”生産者”等に対して、それぞれの属性値を得る。

[0069] 次に、情報処理サーバ10は、属性・属性値に基づき商品等カタログを生成する（ステップS 7）。具体的には、情報処理サーバ10のシステム制御部14は、図10に示すように、商品毎に属性名を並べ、属性名と属性値を組にして商品等カタログを生成する。なお、図10に示すように、商品の画像を商品等カタログに加えてもよい。属性名の順番は、後述する属性のスコアに基づき決定してもよい。

[0070] このように情報処理サーバ10のシステム制御部14は、抽出された属性関連語に基づき、属性関連語が記載された商品等カタログを生成するカタログ生成手段の一例として機能する。

[0071] 次に、情報処理サーバ10のシステム制御部14は、他のカテゴリのウェブページに対しても、ステップS 1～ステップS 7を適用して、商品等カタログを生成する。そして、情報処理サーバ10のシステム制御部14は、生成した商品等カタログの情報を、情報提供サーバ20に送信し、商品等カタログDB22cに記憶させる。

[0072] （3. 2. 属性の選定）

次に、属性の選定のサブルーチンについて、図11を用いて説明する。

[0073] 図11に示すように、情報処理サーバ10は、属性へのスコア付けを行う（ステップS 10）。具体的には、ショッピングサイトが商品を販売する複数の店舗を有する場合、すなわち、サイバーモールを構成する場合、情報処

理サーバ10のシステム制御部14は、属性名が出現したウェブページを有する店舗の数を求め、属性のスコアとする。

[0074] 多種の店舗のウェブページに出現した属性関連語の一例の属性名は、属性として適切であるという仮定に基づいている。例えば、ワインのウェブページにおいて、適切な属性である”品種”という属性は多種の店舗のウェブページに出現する。それに対して、いずれかの属性記述ターンにマッチした不適切な属性は、1店舗のウェブページからしか獲得されないことが多く、属性のスコアが低くなる傾向がある。このように情報処理サーバ10のシステム制御部14は、属性関連語のスコア付けを行う属性スコアリング手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、ウェブページに記載されている対象を販売する複数の店舗を有するウェブサイトにおいて、属性関連語が出現しているウェブページの店舗の数に基づき、属性関連語のスコア付けを行う属性スコアリング手段の一例として機能する。

[0075] 次に、情報処理サーバ10は、上位ランクの属性を選択する（ステップS11）。具体的には、情報処理サーバ10のシステム制御部14は、属性のスコアの高い順に属性名をランク付けし、所定のランク以上の属性名を選択する。このように情報処理サーバ10のシステム制御部14は、スコアの順に属性関連語のランク付けを行い、所定のランク以上の属性関連語を選択する属性選択手段の一例として機能する。

[0076] 次に、情報処理サーバ10は、属性のフィルタリングを行う（ステップS12）。具体的には、情報処理サーバ10のシステム制御部14は、各カテゴリにおける属性名の出現確率を用いて、属性のフィルタリングを行う。他のカテゴリにおいても出現する属性名は、属性として不向きであるという仮定に基づいて、属性のフィルタリングが行われている。例えば、属性として不向きな”送料無料”のような語句は、多数のウェブページに出現するため、各カテゴリにおける出現確率が、似通った値になる。一方、”品種”という属性名はワインのカテゴリのウェブページにはよく出現するが、ゴルフド

ライバーや靴等のカテゴリには出現しないため、ワインのカテゴリにおける出現確率が、ワイン以外のカテゴリにおける出現確率よりも高くなる。このように情報処理サーバ10のシステム制御部14は、カテゴリとは異なるカテゴリに属している複数のウェブページにおいて出現する属性関連語を取り除く属性フィルタ手段の一例として機能する。

[0077] 次に、情報処理サーバ10は、同義の属性を集約する（ステップS13）。属性の中には同じ概念を持つものが存在している。例えば、ワインのカテゴリにおいて、“品種”、“ぶどう品種”、“ブドウ品種”、“セパージュ”、“葡萄品種”は同義の属性名である。情報処理サーバ10のシステム制御部14は、同義語辞書を用いたり、属性名同士の類似の度合いを算出したり、属性名に対応する属性値を用いたりして、同義の属性の属性名を集約する。なお、同義の属性の属性名を集約するのではなく、類似概念の属性の属性名を集約してもよい。

[0078] 具体的には、属性名“A”（属性A）と属性名“B”（属性B）との類似の度合いを算出する場合、属性Aの属性値の中で、属性Bが持っている属性値と共通なものの割合と、属性Bの属性値の中で属性Aの属性値が持っている属性値と共通なものの割合を掛け合わせた値を類似の度合いとしたり、これらの割合を元にエントロピーを計算して掛け合わせた値を類似の度合いとしたり、ジャカード係数を類似の度合いとしたり、属性Aと属性Bの属性値中で共通なものの種類の数を類似の度合いとする。

[0079] このように情報処理サーバ10のシステム制御部14は、属性名同士が類似であるか否かを判定する属性名類似判定手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、属性名類似判定手段により類似と判定された属性名を集約する属性名集約手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、属性関連語として、属性名および属性名に対応する属性値を抽出する属性抽出手段、および、属性値に基づき属性名を集約する属性名集約手段の一例として機能する。

[0080] 本実施形態によれば、ウェブページに記載されている対象を分類するカテ

ゴリが同一である複数のウェブページを取得し、初期データとして、ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得し、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出し、複数のウェブページから、属性関連語に適合する属性記述パターンを抽出することにより、同一の前記カテゴリに属している複数のウェブページから、属性関連語を抽出し、属性記述パターンを抽出するか、または、属性記述パターンを抽出し、属性関連語を抽出しているので、同一のカテゴリに含まれる属性を精度良く取得できる。例えば、属性関連語および属性記述パターンを相互に繰り返し抽出すると、同一のカテゴリに含まれる属性を精度良く取得できる。

- [0081] 情報処理サーバ10のシステム制御部14が、属性抽出手段および属性記述パターン抽出手段を交互に繰り返す場合、属性リストやパターンリストをブートストラップによって拡張して、初期値として与えた属性以外の属性を抽出することができる。また、この抽出された属性により、ウェブページの類似度が判定できる。また、ユーザがウェブページに関する商品等カタログを使用して、所望の商品に到達しやすくなり、ユーザの利便性の向上を図ることができる。
- [0082] また、情報処理サーバ10のシステム制御部14が、抽出された属性関連語から属性リストを生成し、抽出された属性記述パターンのパターンリストを生成する場合、カテゴリ毎に、属性名や属性値等の属性関連語や属性記述パターンの情報を蓄積できる。
- [0083] また、情報処理サーバ10のシステム制御部14が、属性関連語のスコア付けを行い、上位のランクの属性関連語を選択する場合、選択された属性関連語において、商品等を表す属性やブログの属性の精度が高くなる。
- [0084] また、情報処理サーバ10のシステム制御部14が、対象を販売する複数の店舗を有するウェブサイトにおいて、属性関連語が出現しているウェブページの店舗の数に基づいて属性関連語のスコア付けを行う場合、属性関連語

を選択する際、商品等を表す属性の精度が高くなる。例えば、店舗により扱う商品等の数や、ウェブページの数が大きく異なる場合、多くの商品等を扱う店舗の影響を受けやすくなるが、店舗の数に基づき属性関連語のスコア付けを行うことにより、ある特定の店舗の影響を解消できる。

[0085] また、情報処理サーバ10のシステム制御部14が、他のカテゴリに属している複数のウェブページにおいて出現する属性関連語を取り除く場合、対象のカテゴリ固有の属性関連語に絞ることにより、商品等を表す属性やブログの属性の精度が高くなる。

[0086] また、情報処理サーバ10のシステム制御部14が、属性関連語として属性の属性名を抽出する場合、同一のカテゴリに含まれる属性・属性名を精度良く取得できる。また、属性名により、誤ったカテゴリに登録されたウェブページを抽出できる。

[0087] また、情報処理サーバ10のシステム制御部14が、属性名同士が類似であるか否かを判定し、類似と判定された属性名を集約する場合、重複した属性名を取り除かれ、属性名が利用しやすくなる。

[0088] また、情報処理サーバ10のシステム制御部14が、属性関連語として、属性名および属性名に対応する属性値を抽出し、属性値に基づき属性名を集約する場合、属性名に直結した属性値により、属性名が集約しやすくなる。

[0089] また、情報処理サーバ10のシステム制御部14が、対象の供給元のウェブページを取得し、初期データ取得手段、属性抽出手段、および、属性記述パターン抽出手段により、対象の供給元のウェブページから供給元対象属性関連語を抽出し、供給元対象属性関連語と属性関連語とを比較する場合、同一のカテゴリに含まれる属性をより精度良く取得できる。また、商品等の対象に関する公式な対象情報を取り入れ、生成されたカタログの精度を判定することにより、カタログの信頼性を向上させることができる。

[0090] また、情報処理サーバ10のシステム制御部14が、抽出された属性関連語に基づき、属性関連語が記載されたカタログを生成する場合、ユーザがウェブページに関するカタログを使用して、所望の商品等の対象に到達しやす

くなり、ユーザの利便性の向上を図ることができる。

[0091] [4. 商品等のカタログ生成システムの動作の第1変形例]

次に、商品等のカタログ生成システムの動作の第1変形例について図12および図13に基づきについて説明する。

[0092] なお、上記実施形態と同一または対応する部分には、同一の符号を用いて動作等を説明する。その他の変形例も同様とする。

[0093] 図12は、商品等のカタログ生成の第1変形例の動作例を示すフローチャートである。図13は、第1変形例の属性・属性値の抽出の様子の一例を示す模式図である。図12に示すように、本変形例は、ブートストラップ法における初期データを、属性記述パターンではなく、属性・属性値とした点である。ステップS22からステップS24が、上記実施形態と異なるステップである。なお、属性・属性値DB12bには、属性・属性値の初期データが記憶されている。

[0094] まず、情報処理サーバ10は、ステップS1と同様に、複数のウェブページを取得する（ステップS21）。

[0095] 次に、情報処理サーバ10は、属性・属性値を取得する（ステップS22）。具体的には、情報処理サーバ10のシステム制御部14は、下記のステップS23～S25におけるブートストラップ法の初期データとして、属性・属性値DB12bの属性・属性値リストから、図13に示すように、初期の属性・属性値66を取得する。このように情報処理サーバ10のシステム制御部14は、初期データとして、商品の属性に関連した属性関連語を取得する初期データ取得手段の一例として機能する。

[0096] 次に、情報処理サーバ10は、属性・属性値に適合する属性記述パターンを抽出する（ステップS23）。具体的には、情報処理サーバ10のシステム制御部14は、ウェブページ50等の複数のウェブページの中から、図13に示すように、属性・属性値66等に適合する語句の部分（例えば”[品種：○○○]”）を取り出し、属性記述パターン”[:]”等を抽出する。そして、抽出した属性記述パターンは、属性記述パターンリストとして

属性記述パターンDB 12 aに記憶される。ここで、ワイルドカードと属性・属性値とが用いられて、属性記述パターンが抽出される。

[0097] 次に、情報処理サーバ10は、属性記述パターンに適合する属性・属性値を抽出する（ステップS 24）。具体的には、情報処理サーバ10のシステム制御部14は、図13に示すように、属性記述パターン67（例えば、属性記述パターンの前部” [”、中部” :”、後部”] ”）に適合する、例えば、” [アルコール度数：12.5%”] を取り出し、属性・属性値をウェブページ50等の複数のウェブページの中から抽出する。そして、抽出した属性・属性値は、属性・属性値リストに追加され、属性・属性値DB 12 bに記憶される。

[0098] 以下のステップS 25からステップS 28は、ステップS 5からステップS 8と同様である。

[0099] 以上のように、情報処理サーバ10のシステム制御部14は、ステップS 22からステップS 24において、商品进行分类するカテゴリにおいて、同一のカテゴリに属している複数のウェブページを取得するウェブページ取得ステップと、商品の属性に関連した属性関連語を取得する属性関連語取得ステップと、属性の記述に用いられる属性記述パターンであって、複数のウェブページから、属性関連語に適合する属性記述パターンを抽出する属性記述パターン抽出ステップと、抽出された属性関連語に基づき、属性記述パターン抽出手段に使用する属性関連語を、複数のウェブページから更に抽出する属性関連語抽出ステップとを実行する。

[0100] 本変形例によれば、商品进行分类するカテゴリにおいて、同一のカテゴリに属している複数のウェブページを取得し、属性・属性値DB 12 bから商品の属性に関連した属性関連語を取得し、属性の記述に用いられる属性記述パターンであって、複数のウェブページから属性関連語に適合する属性記述パターンを抽出し、抽出された属性関連語に基づき、属性記述パターンの抽出に使用する属性関連語を、複数のウェブページから更に抽出し、抽出された属性関連語に基づき、属性関連語が記載された商品等カタログを生成するこ

とにより、ユーザがウェブページに関する商品等カタログを使用して、所望の商品に到達しやすくなり、ユーザの利便性の向上を図ることができる。

[0101] [5. 商品等のカタログ生成システムの動作の第2変形例]

次に、商品等のカタログ生成システムの動作の第2変形例について図14～図16に基づきについて説明する。本変形例では、ブートストラップ法のステップにおいて、属性の選定を行ったり、属性記述パターンの選定を行ったりしている。

[0102] 図14は、商品等のカタログ生成の第2変形例の動作例を示すフローチャートである。図15は、属性・属性値抽出のサブルーチンの一例を示すフローチャートである。図16は、属性記述パターン抽出のサブルーチンの一例を示すフローチャートである。

[0103] (5. 1. ウェブページの抽出の流れ)

まず、図14に示すように、情報処理サーバ10は、ステップS1およびステップS2と同様に、複数のウェブページを取得し（ステップS31）、属性記述パターンを取得する（ステップS32）。

[0104] 次に、情報処理サーバ10は、属性記述パターンに基づき属性・属性値を抽出する（ステップS33）。具体的には、情報処理サーバ10のシステム制御部14は、属性・属性値抽出のサブルーチンにより属性・属性値を抽出する。属性・属性値抽出のサブルーチンでは、情報処理サーバ10のシステム制御部14は、属性記述パターンに適合する属性・属性値を抽出したり、属性へのスコア付けを行ったり、上位のランクの属性を選択したり、属性のフィルタリングを行ったりする。

[0105] 次に、情報処理サーバ10は、属性・属性値に基づき属性記述パターンを抽出する（ステップS34）。具体的には、情報処理サーバ10のシステム制御部14は、属性記述パターン抽出のサブルーチンにより属性記述パターンを抽出する。属性記述パターン抽出のサブルーチンでは、情報処理サーバ10のシステム制御部14は、属性・属性値に適合する属性記述パターンを抽出したり、属性記述パターンと属性・属性値との共起確率を算出したり、

スコアを算出したり、上位のランクの属性記述パターンを選択したりする。

[0106] 次に、情報処理サーバ10は、ステップS5と同様に、所定回数を判定する（ステップS35）。

[0107] 次に、情報処理サーバ10は、同義の属性を集約する（ステップS36）。具体的には、情報処理サーバ10のシステム制御部14は、ステップS33～ステップS35のブートストラップ法により求めた属性名に対して、ステップS13と同様に、同義の属性の属性名の集約を行う。

[0108] 次に、情報処理サーバ10は、ステップS7と同様に、属性、属性値に基づき、商品等カタログを生成する（ステップS37）。

[0109] （5. 2. 属性・属性値の抽出）

次に、属性・属性値抽出のサブルーチンについて、図15を用いて説明する。

[0110] 図15に示すように、情報処理サーバ10は、属性記述パターンに適合する属性・属性値を抽出する（ステップS40）。具体的には、情報処理サーバ10のシステム制御部14は、ステップS3と同様に、属性記述パターンに適合する属性・属性値を抽出する。

[0111] 次に、情報処理サーバ10は、属性選定のサブルーチンにおけるステップS10～ステップS12と同様に、属性へのスコア付けを行い（ステップS41）、上位のランクの属性を選択し（ステップS42）、属性のフィルタリングを行う（ステップS43）。

[0112] （5. 3. 属性記述パターンの抽出）

次に、属性記述パターン抽出のサブルーチンについて、図16を用いて説明する。

[0113] 図16に示すように、情報処理サーバ10は、ステップS4と同様に属性・属性値に適合する属性記述パターンを抽出する（ステップS45）。

[0114] 次に、情報処理サーバ10は、属性記述パターンと属性・属性値との共起確率を算出する（ステップS46）。具体的には、情報処理サーバ10のシステム制御部14は、属性関連語と属性記述パターンとが共に出現する共起

数の一例として、属性記述パターン t と、対象となっている同一カテゴリの複数のウェブページにおいて、属性・属性値の対 i との共起数 $f(i, t)$ を算出する。そして、情報処理サーバ 10 のシステム制御部 14 は、属性記述パターン t と属性・属性値の組 i が共起する確率、すなわち、式 (1) の共起確率 $P_t(i)$ を算出する。

$$P_t(i) = f(i, t) / N_t \quad \dots (1)$$

ここで、 N_t は、抽出した属性記述パターン t の数である。

[0115] 次に、情報処理サーバ 10 は、エントロピー（スコア）を算出する（ステップ S 47）。様々な属性・属性値と共起する属性記述パターンは、適切な属性記述パターンであるという仮定に基づいて、情報処理サーバ 10 のシステム制御部 14 は、属性記述パターンに対するエントロピー $H(t)$ を式 (2) により算出する。

$$H(t) = - \sum_{i \in I} P_t(i) \times \log_2 P_t(i) \quad \dots (2)$$

[0116] ここで、 I は、属性・属性値の組 i を要素とする属性・属性値の組の集合、属性記述パターン t を要素とする T は属性記述パターン集合である。

[0117] 次に、情報処理サーバ 10 は、上位のスコアの属性記述パターンを選択する（ステップ S 48）。具体的には、情報処理サーバ 10 のシステム制御部 14 は、スコアとしてエントロピー $H(t)$ の高い属性記述パターンからランク付けをして、所定のランク以上の属性記述パターンを選択する。このように情報処理サーバ 10 のシステム制御部 14 は、属性記述パターンのスコア付けを行う属性記述パターン・スコアリング手段の一例として機能する。また、情報処理サーバ 10 のシステム制御部 14 は、スコアの順に属性記述パターンのランク付けを行い、所定のランク以上の属性記述パターンを選択する属性記述パターン選択手段の一例として機能する。また、情報処理サーバ 10 のシステム制御部 14 は、属性関連語と属性記述パターンとが共に出現する共起数に基づき属性記述パターンのスコア付けを行う属性記述パターン・スコアリング手段の一例として機能する。

[0118] なお、ステップ S 46～ステップ S 48 は、ステップ S 5 までに得られた

属性記述パターンリストに対して、属性記述パターンの選定のステップとして使用されて、無駄な属性記述パターンを省くことができる。

- [0119] 以上のように、本変形例では、特に、ブートストラップの回数を増やした場合に、抽出される属性・属性値や、属性記述パターンが増大することを防止することができる。
- [0120] また、情報処理サーバ10のシステム制御部14が、属性記述パターンのスコア付けを行い、上位のランクの属性記述パターンを選択する場合、属性・属性値を抽出するための属性記述パターンの精度が高くなる。
- [0121] また、情報処理サーバ10のシステム制御部14が、属性関連語と属性記述パターンとが共に出現する共起数に基づき属性記述パターンのスコア付けを行う場合、属性記述パターンを選択する上でのスコアの精度が高くなる。
- [0122] なお、本実施形態やその変形例において、属性のスコアとして、店舗数でなく、属性名が出現したウェブページの数でもよい。情報処理サーバ10のシステム制御部14が、属性スコアリング手段として、属性関連語の検索のヒット件数に基づき、属性関連語のスコア付けを行う。この場合、店舗が多く集まるサイバーモール以外にも適用できる。
- [0123] また、商品等の対象の供給元のウェブページから、各商品等の対象の属性・属性値を求め、この属性・属性値により、商品等カタログの属性・属性値の精度の判定を行ってもよい。この場合、商品等に関する公式な商品等情報を取り入れ、生成された商品等カタログの精度を判定することにより、商品等カタログの信頼性を向上させることができる。
- [0124] 例えば、図17に示すように、情報処理サーバ10のが、商品等の供給元対象の製造元や輸入販売元等の対象供給元のウェブページを通信部を通して取得する（ステップS51）。具体的には、情報処理サーバ10のシステム制御部14が、情報DB22aに記憶されたURL等を参照して、図18に示すような、各商品等の商品IDに対応した供給元のウェブページを取得する。なお、商品IDは、ショッピングサイトにおいて、ウェブページに予め付されている商品IDや、ウェブページのテキストデータから抽出した商品

IDでもよい。このように、情報処理サーバ10のシステム制御部14および通信部11が、商品等の対象の供給元のウェブページを取得するウェブページ取得手段の一例として機能する。

[0125] 次に、情報処理サーバ10のが、供給元商品の属性・属性値を抽出する（ステップS52）。具体的には、情報処理サーバ10のシステム制御部14が、ステップS2～ステップS6や、ステップS22～ステップS26や、ステップS32～ステップS36のようにして、商品供給元のウェブページから、商品供給元のウェブページに関する属性・属性値を抽出する。

[0126] 次に、情報処理サーバ10のが、供給元商品等の供給元対象の属性・属性値と商品等カタログの属性・属性値とを比較して、商品等カタログの精度を判定する。具体的には、情報処理サーバ10のシステム制御部14が、商品等カタログの属性名に、供給元商品等の供給元対象の属性名があるか否か、また、属性名がある場合、その属性値が一致しているか否かを比較する。そして、同一の属性名が少ない場合や属性値が一致しない場合は、生成された商品等カタログの精度が低いと判断する。また、属性値が一致しない場合は、店舗側での入力ミスと考えることができる。そして、属性名が一致している数や、属性値が一致している数に所定値を設け、所定値以上の場合の商品等カタログの情報が、ユーザ端末35から閲覧されるようにする。このように、情報処理サーバ10のシステム制御部14および通信部11が、初期データ取得手段、属性抽出手段、および、属性記述パターン抽出手段により、商品等の対象の供給元のウェブページから供給元対象属性関連語を抽出し、供給元対象属性関連語と、属性関連語とを比較する属性関連語比較手段の一例として機能する。

[0127] 次に、生成された商品等カタログの変形例について図19を用いて説明する。

図19は、生成された商品等カタログの一例を示す模式図である。

[0128] 図19に示すように、属性名”製造年”に対して、属性値”1995年”および属性値”1996年”によりグルーピングされている。情報処理サー

サーバ10のシステム制御部14が、ステップS6や、ステップS26や、ステップS36等で、属性関連語を求めた後、属性名”製造年”対して属性値”1995年”を有するウェブページを収集し、また、属性名”製造年”対して属性値”1996年”を有するウェブページを収集して、属性関連語に基づき、複数のウェブページの情報グルーピングする。図19に示すように、商品の商品名と、他の属性とが、属性名”製造年”対する属性値によりグルーピングされる。

[0129] 情報処理サーバ10のシステム制御部14が、属性関連語に基づき、複数のウェブページをグルーピングする場合、共通する属性によりグルーピングされたウェブページにまとめることができる。また、ユーザが見やすい検索結果に反映させる等、利用価値が高まる。

[0130] (第2実施形態)

次に、本発明の第2実施形態に係る情報処理システム1の動作について、図を用いて説明する。なお、前記第1実施形態と同一または対応する部分には、同一の符号を用いて異なる構成および作用のみを説明する。その他の実施形態および変形例も同様とする。

[0131] [6. 第2実施形態に係るウェブページ抽出システムの動作]

次に、本発明の第2実施形態に係る情報処理システム1のウェブページ抽出システムとしての動作について図20を用いて説明する。

[0132] まず、ウェブページの抽出の流れについて、図20を用いて説明する。

図20は、第2実施形態に係る情報処理システムにおいてウェブページを抽出する動作例を示すフローチャートである。

[0133] 図20に示すように、情報処理サーバ10は、複数のウェブページを取得し、属性を求める(ステップS60)。具体的には、情報処理サーバ10のシステム制御部14は、ステップS1～ステップS6と同様に複数のウェブページを取得し、選定した属性の属性名および属性値を求める。または、情報処理サーバ10のシステム制御部14は、ステップS21～ステップS26と同様に複数のウェブページを取得し、選定した属性の属性名および属性

値を求める。

- [0134] 次に、情報処理サーバ10は、各ウェブページの属性の出現回数を算出する（ステップS61）。具体的には、情報処理サーバ10のシステム制御部14は、ステップS60で取得した各ウェブページで、選定した属性の属性名の出現回数を算出する。なお、集約された属性名の同義語も考慮して、出現回数が算出される。
- [0135] 次に、情報処理サーバ10は、出現回数が所定回数以下のウェブページを抽出する（ステップS62）。具体的には、情報処理サーバ10のシステム制御部14は、ステップS60で取得したウェブページで、選定した属性の属性名の出現回数がゼロ、すなわち、選定した属性の属性名が出現しないウェブページを抽出する。複数の属性名がある場合は、どの属性名も出現しないウェブページを抽出する。情報処理サーバ10のシステム制御部14は、この抽出されたウェブページは誤ったカテゴリに登録されたウェブページである、と判定する。
- [0136] このように情報処理サーバ10のシステム制御部14は、複数のウェブページから、属性関連語の出現回数が所定回数以下のウェブページを抽出するウェブページ抽出手段の一例として機能する。また、情報処理サーバ10のシステム制御部14は、属性関連語の出現回数がゼロのウェブページを抽出するウェブページ抽出手段の一例として機能する。
- [0137] 次に、情報処理サーバ10のシステム制御部14は、他のカテゴリのウェブページに対しても、ステップS60～ステップS62を適用して、誤ったカテゴリに登録されたウェブページを抽出する。そして、情報処理サーバ10のシステム制御部14は、抽出したウェブページに関する情報を、店舗主に知らせるために店舗主端末30に送信したり、抽出したウェブページを正しいカテゴリに移動させたりする。
- [0138] 本実施形態によれば、ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得し、初期データとして、ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対

象の属性の記述に用いられる属性記述パターンを取得し、複数のウェブページから、属性記述パターンに適合する属性の属性関連語を抽出し、複数のウェブページから、属性関連語に適合する前記属性記述パターンを抽出し、複数のウェブページから、属性関連語の出現回数が所定回数以下のウェブページを抽出することにより、誤ったカテゴリに登録されたウェブページを抽出することができる。例えば、ワインセラーの商品がワインのカテゴリに登録された場合、ワインセラーに関するウェブページには、ワインの属性のひとつ”品種”等の属性名が出てこない確率が高い。また、属性の属性名を抽出する場合、属性名により、誤ったカテゴリに登録された商品ウェブページを抽出できる。

[0139] また、情報処理サーバ10のシステム制御部14が、属性関連語の出現回数がゼロのウェブページを抽出する場合、誤ったカテゴリに登録されたウェブページには、属性関連語が出現する確率が低いので、誤ったカテゴリに登録されたウェブページを容易に抽出することができる。

[0140] [7. ウェブページ抽出システムの動作の第1変形例]

次に、ウェブページ抽出システムの動作の第1変形例について図21を用いて説明する。本変形例では、ブートストラップ法のステップにおいて、属性の選定を行ったり、属性記述パターンの選定を行ったりしている。

[0141] 図21は、ウェブページ抽出の第1変形例の動作例を示すフローチャートである。

[0142] (5. 1. ウェブページの抽出の流れ)

まず、図21に示すように、情報処理サーバ10は、ステップS31～ステップS36と同様に、複数のウェブページを取得し、属性を求め、同義の属性を集約する(ステップS65)。

[0143] 次に、情報処理サーバ10は、ステップS61およびステップS62と同様に、各ウェブページの属性の属性名の出現回数を算出し(ステップS66)、出現回数が所定回数以下のウェブページを抽出する(ステップS67)。

- [0144] 以上のように、本変形例では、特に、ブートストラップの回数を増やした場合に、抽出される属性・属性値や、属性記述パターンが増大することを防止することができる。
- [0145] また、情報処理サーバ10のシステム制御部14が、属性記述パターンのスコア付けを行い、上位のランクの属性記述パターンを選択する場合、属性・属性値を抽出するための属性記述パターンの精度が高くなる。
- [0146] また、情報処理サーバ10のシステム制御部14が、属性関連語と属性記述パターンとの共起数に基づき属性記述パターンのスコア付けを行う場合、属性記述パターンを選択する上でのスコアの精度が高くなる。
- [0147] なお、本実施形態やその変形例において、属性のスコアとして、店舗数でなく、属性名が出現したウェブページの数でもよい。情報処理サーバ10のシステム制御部14が、属性スコアリング手段として、属性関連語の検索のヒット件数に基づき、属性関連語のスコア付けを行う。この場合、店舗が多く集まるサイバーモール以外にも適用できる。
- [0148] また、情報処理サーバ10のシステム制御部14が、ウェブページ抽出手段として、複数の属性関連語において、属性関連語の出現回数がゼロの割合が、所定以上のウェブページを抽出してもよい。また、属性関連語の出現回数がゼロでなくても、少数出現回数に基づいてもよい。
- [0149] さらに、本発明は、上記各実施形態に限定されるものではない。上記各実施形態は、例示であり、本発明の特許請求の範囲に記載された技術的思想と実質的に同一な構成を有し、同様な作用効果を奏するものは、いかなるものであっても本発明の技術的範囲に包含される。

符号の説明

- [0150] 3 : ネットワーク
5 : サーバシステム
10 : 情報処理サーバ (情報処理装置)
12 a : 属性記述パターンDB
12 b : 属性・属性名DB

20 : 情報提供サーバ

22c : 商品等カタログDB

請求の範囲

- [請求項1] ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段と、
初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段と、
前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段と、
前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段と、
を備えたことを特徴とする情報処理装置。
- [請求項2] 請求項1に記載の情報処理装置において、
前記属性抽出手段および前記属性記述パターン抽出手段を交互に繰り返す繰返手段を更に備えたことを特徴とする情報処理装置。
- [請求項3] 請求項1または請求項2に記載の情報処理装置において、
前記属性抽出手段が、前記属性関連語として、前記属性の属性名を抽出することを特徴とする情報処理装置。
- [請求項4] 請求項1から請求項3のいずれか1項に記載の情報処理装置において、
抽出された前記属性関連語から属性リストを生成する属性リスト生成手段と、
抽出された前記属性記述パターンのパターンリストを生成するパターンリスト生成手段と、
を更に備えたことを特徴とする情報処理装置。
- [請求項5] 請求項1から請求項4のいずれか1項に記載の情報処理装置において、
前記属性関連語のスコア付けを行う属性スコアリング手段と、
前記スコアの順に前記属性関連語のランク付けを行い、所定のランク

以上の属性関連語を選択する属性選択手段と、
を更に備えたことを特徴とする情報処理装置。

[請求項6] 請求項5に記載の情報処理装置において、
前記属性スコアリング手段が、前記属性関連語の検索のヒット件数に基づき、前記属性関連語のスコア付けを行うことを特徴とする情報処理装置。

[請求項7] 請求項5に記載の情報処理装置において、
前記属性スコアリング手段が、前記対象を販売する複数の店舗を有するウェブサイトにおいて、前記属性関連語が出現しているウェブページの前記店舗の数に基づき、前記属性関連語のスコア付けを行うことを特徴とする情報処理装置。

[請求項8] 請求項1から請求項7のいずれか1項に記載の情報処理装置において、
前記カテゴリとは異なるカテゴリに属している複数のウェブページにおいて出現する前記属性関連語を取り除く属性フィルタ手段を更に備えたことを特徴とする情報処理装置。

[請求項9] 請求項1から請求項8のいずれか1項に記載の情報処理装置において、
前記属性記述パターンのスコア付けを行う属性記述パターン・スコアリング手段と、
前記スコアの順に前記属性記述パターンのランク付けを行い、所定のランク以上の属性記述パターンを選択する属性記述パターン選択手段と、
を更に備えたことを特徴とする情報処理装置。

[請求項10] 請求項9に記載の情報処理装置において、
前記属性記述パターン・スコアリング手段が、前記属性関連語と前記属性記述パターンとが共に出現する共起数に基づき前記属性記述パターンのスコア付けを行うことを特徴とする情報処理装置。

- [請求項11] 請求項 1 から請求項 10 のいずれか 1 項に記載の情報処理装置において、
- 前記属性名同士が類似であるか否かを判定する属性名類似判定手段と、
- 前記属性名類似判定手段により類似と判定された属性名を集約する属性名集約手段と、
- を更に備えたことを特徴とする情報処理装置。
- [請求項12] 請求項 11 に記載の情報処理装置において、
- 前記属性抽出手段が、前記属性関連語として、前記属性名および前記属性名に対応する属性値を抽出し、
- 前記属性名集約手段が、前記属性値に基づき前記属性名を集約することを特徴とする情報処理装置。
- [請求項13] 請求項 1 から請求項 10 のいずれか 1 項に記載の情報処理装置において、
- 前記ウェブページ取得手段が、前記対象の供給元のウェブページを取得し、
- 前記初期データ取得手段、前記属性抽出手段、および、前記属性記述パターン抽出手段により、前記対象の供給元のウェブページから供給元対象属性関連語を抽出し、前記供給元対象属性関連語と前記属性関連語とを比較する属性関連語比較手段を更に備えたことを特徴とする情報処理装置。
- [請求項14] 請求項 1 から請求項 13 のいずれか 1 項に記載の情報処理装置において、
- 抽出された前記属性関連語に基づき、前記属性関連語が記載されたカタログを生成するカタログ生成手段を更に備えたことを特徴とする情報処理装置。
- [請求項15] 請求項 1 から請求項 14 のいずれか 1 項に記載の情報処理装置において、

前記複数のウェブページから、前記属性関連語の出現回数が所定回数以下のウェブページを抽出するウェブページ抽出手段と、
を更に備えたことを特徴とする情報処理装置。

[請求項16]

請求項15に記載の情報処理装置において、
前記ウェブページ抽出手段が、前記属性関連語の出現回数がゼロのウェブページを抽出することを特徴とする情報処理装置。

[請求項17]

請求項1から請求項16のいずれか1項に記載の情報処理装置において、
前記属性関連語に基づき、前記複数のウェブページをグルーピングするウェブページ・属性グルーピング手段を更に備えたことを特徴とする情報処理装置。

[請求項18]

情報処理装置が情報処理をする情報処理方法において、
ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得ステップと、
前記ウェブページに記載されている対象の属性の記述に用いられる属性記述パターンを取得する属性記述パターン取得ステップと、
前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出ステップと、
抽出された前記属性関連語に基づき、前記属性抽出ステップで使用する前記属性記述パターンを、前記複数のウェブページから、更に抽出する属性記述パターン抽出ステップと、
を有することを特徴とする情報処理方法。

[請求項19]

情報処理装置が情報処理をする情報処理方法において、
ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得ステップと、
前記ウェブページに記載されている対象の属性に関連した属性関連語を取得する属性関連語取得ステップと、
前記属性の記述に用いられる属性記述パターンであって、前記複数の

のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出ステップと、

抽出された前記属性関連語に基づき、前記属性記述パターン抽出ステップで使用する属性関連語を、前記複数のウェブページから、更に抽出する属性関連語抽出ステップと、

を有することを特徴とする情報処理方法。

[請求項20]

コンピュータを、

ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段、

初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段、

前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段、および、

前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段として機能させることを特徴とする情報処理装置用のプログラム。

[請求項21]

コンピュータを、

ウェブページに記載されている対象を分類するカテゴリが同一である複数のウェブページを取得するウェブページ取得手段、

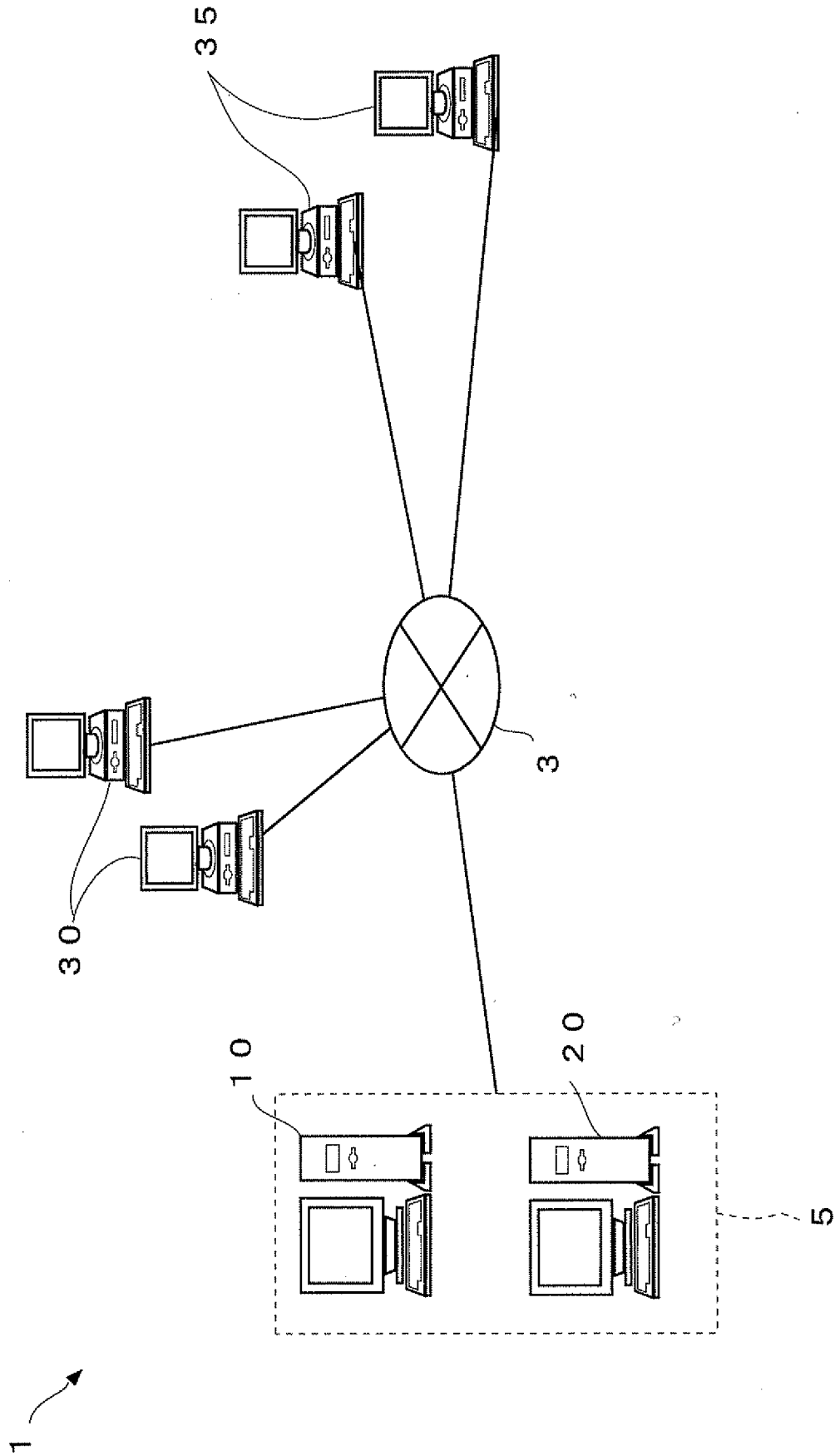
初期データとして、前記ウェブページに記載されている対象の属性に関連した属性関連語、または、当該対象の属性の記述に用いられる属性記述パターンを取得する初期データ取得手段、

前記複数のウェブページから、前記属性記述パターンに適合する前記属性の属性関連語を抽出する属性抽出手段、および、

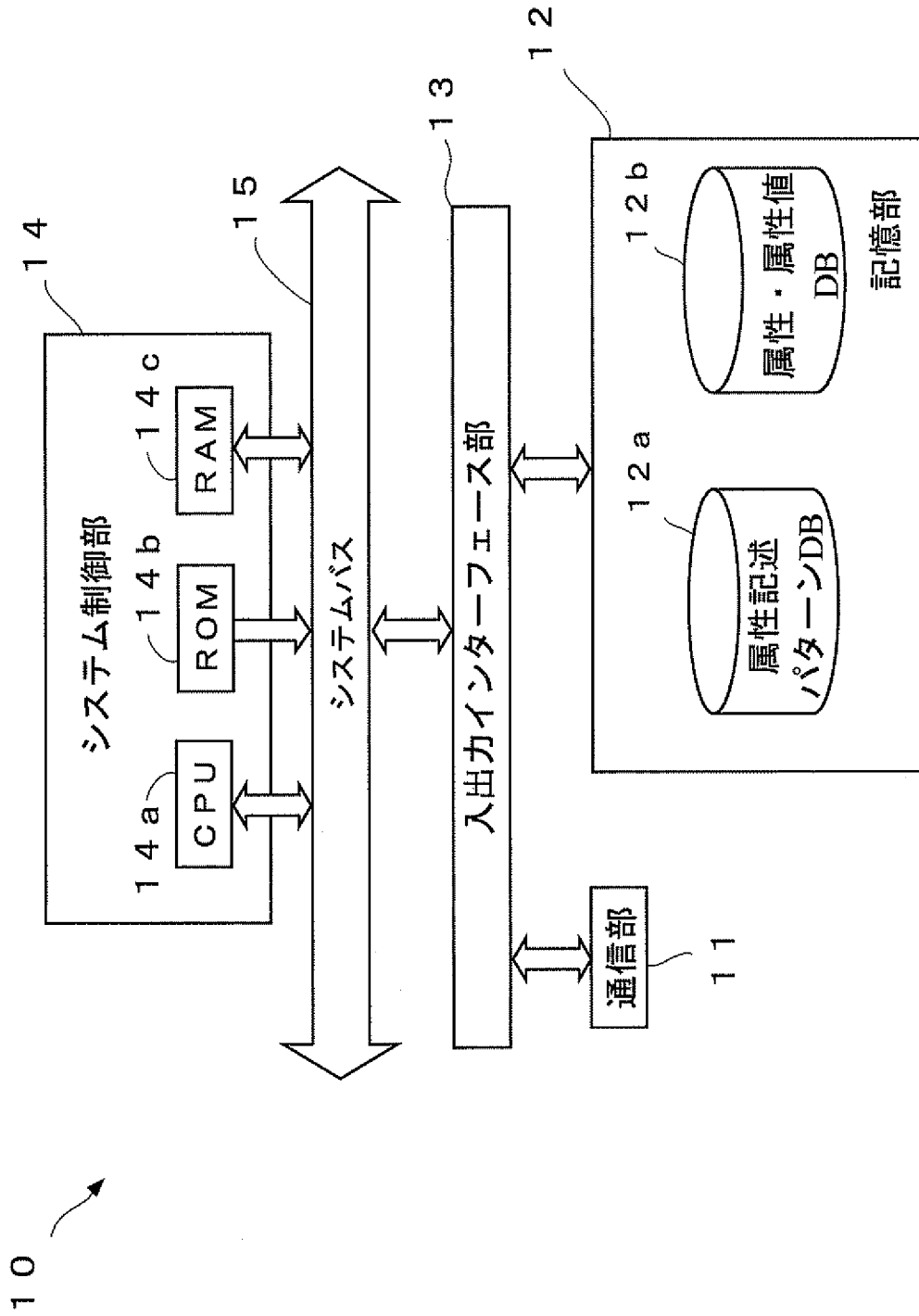
前記複数のウェブページから、前記属性関連語に適合する前記属性記述パターンを抽出する属性記述パターン抽出手段として機能させることを特徴とする情報処理装置用のプログラムを記録したコンピュー

タ読み取り可能な記録媒体。

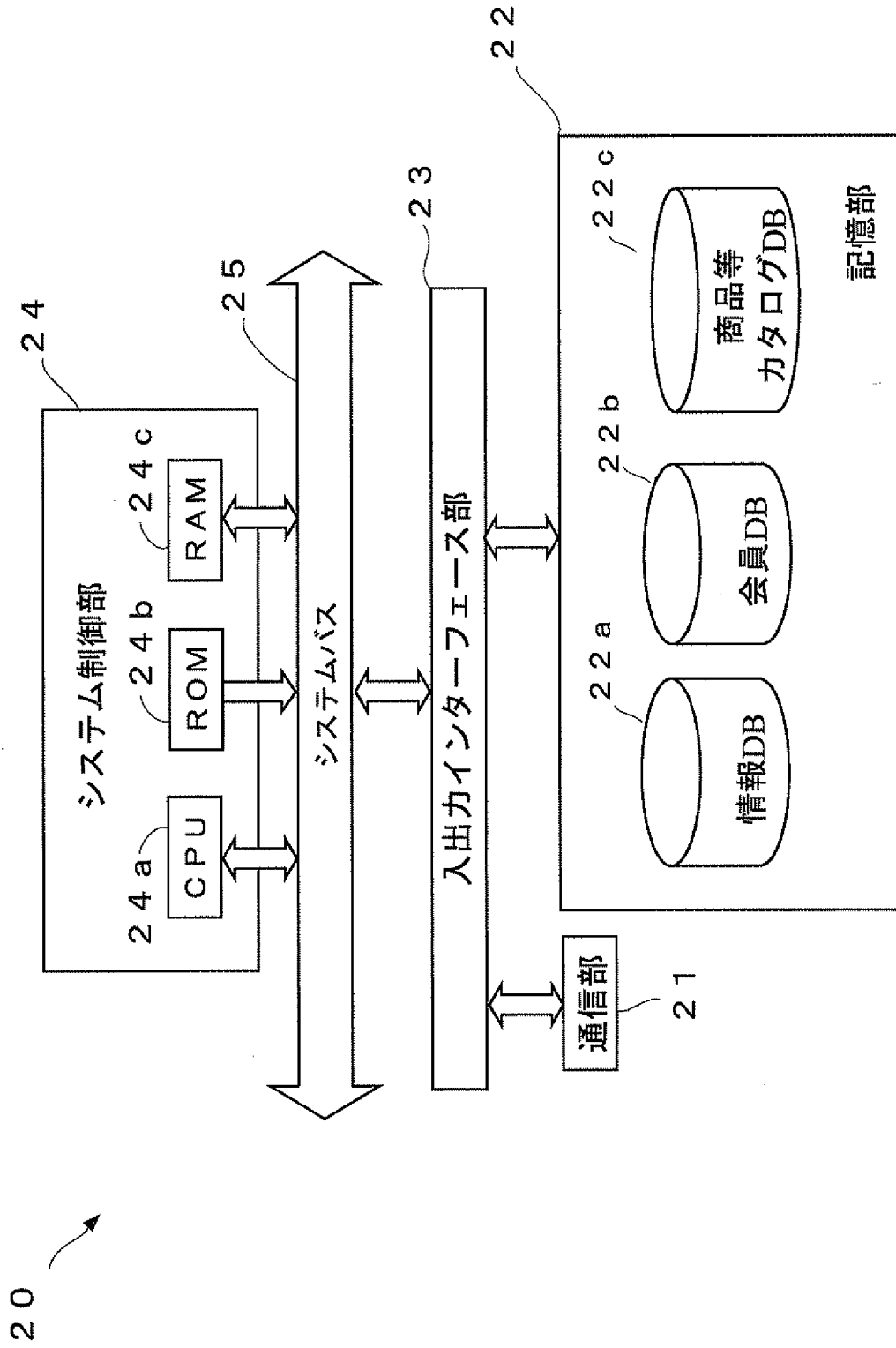
[図1]



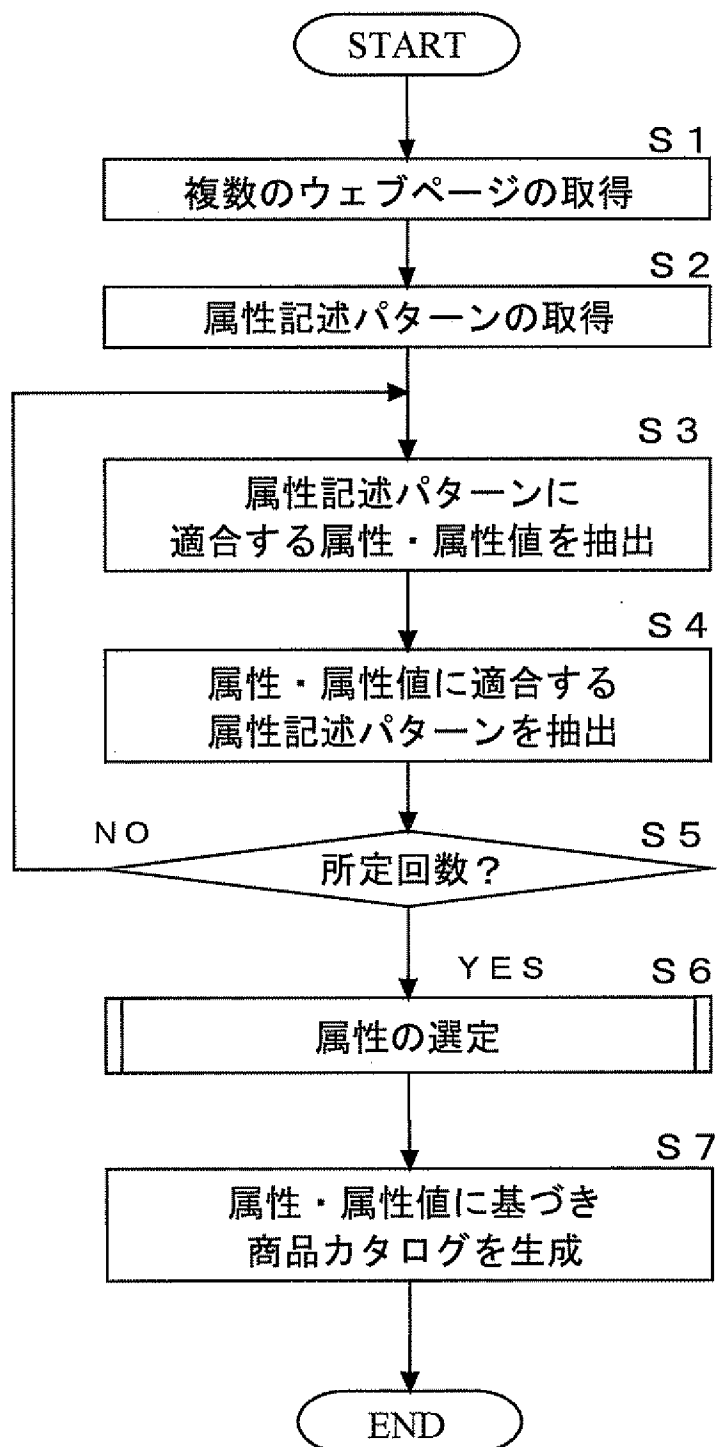
[図2]



[図3]



[図4]



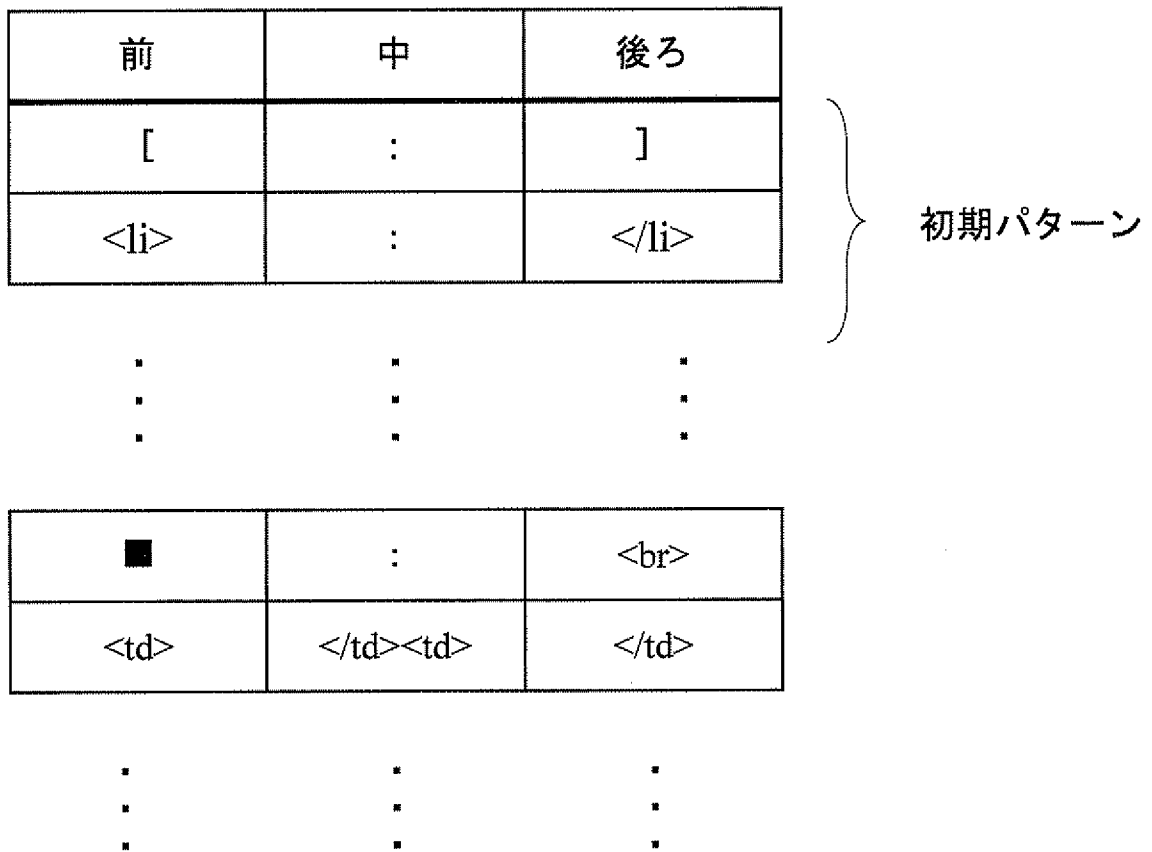
[図5]



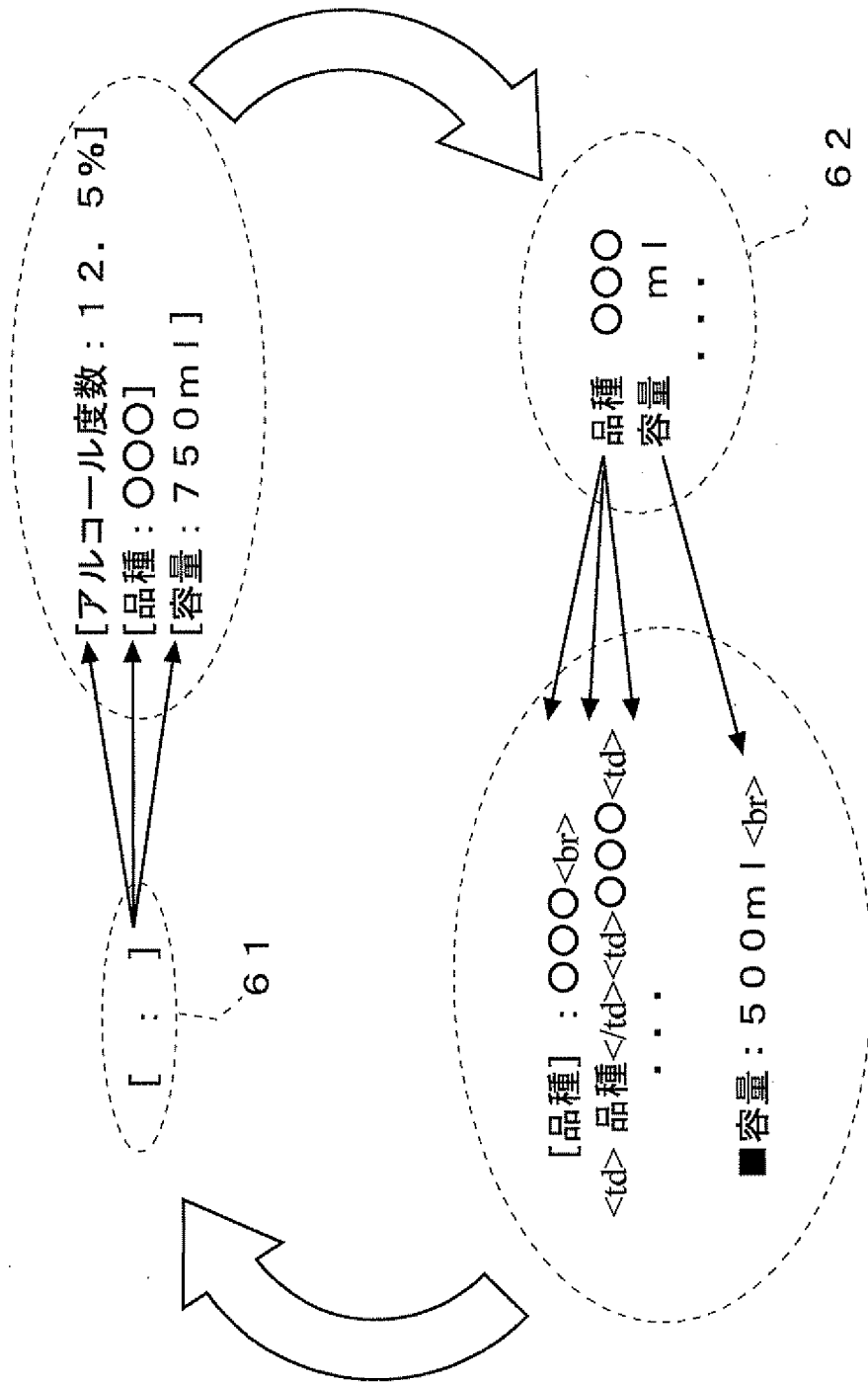
[図6]

```
<html>
<head>
  . . . . .
<title> [〇〇市場] . . . 。 ツエ〇・〇〇・△△ . .
  . . . . .
</title>
  . . . . .
</head>
<body>
  . . . . .
  . . . . ツエ〇・〇〇・△△ . . . . .
  . . . . .
  . . . . . [送料無料] . . . . .
  . . . . . [品種：〇〇〇] . . . . .
  . . . . .
  . . . . .
  . . . . .
<table . . . . .>
  . . . . .
  <td>〇〇</td><td>◇◇</td>
  . . . . .
</table>
  . . . . .
</body>
</html>
```

[図7]



[図8]



[図9]

ワイン

属性名	属性値
品種	〇〇〇、〇▼
生産者	△△、…
アルコール度数	12.5%

⋮

⋮

ゴルフドライバー

属性名	属性値
重量	50g、…
ロフト角	10.5°、…
シャフト	…

⋮

⋮

[図10]



●●ワイン

商品番号1 2 3 4 5

品種	〇〇
生産者	△△
容量	750ml
製造年	2000年
価格	¥912

⋮
⋮
⋮

⋮
⋮
⋮

▽▲ワイン

商品番号1 2 3 4 6

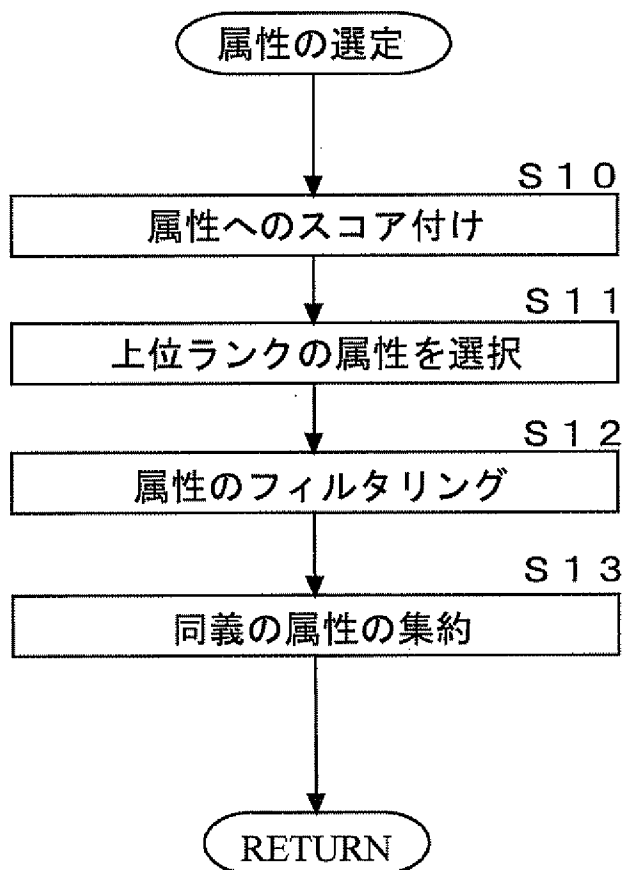


品種	〇●
生産者	▲△
容量	750ml
製造年	1995年
価格	¥2,900

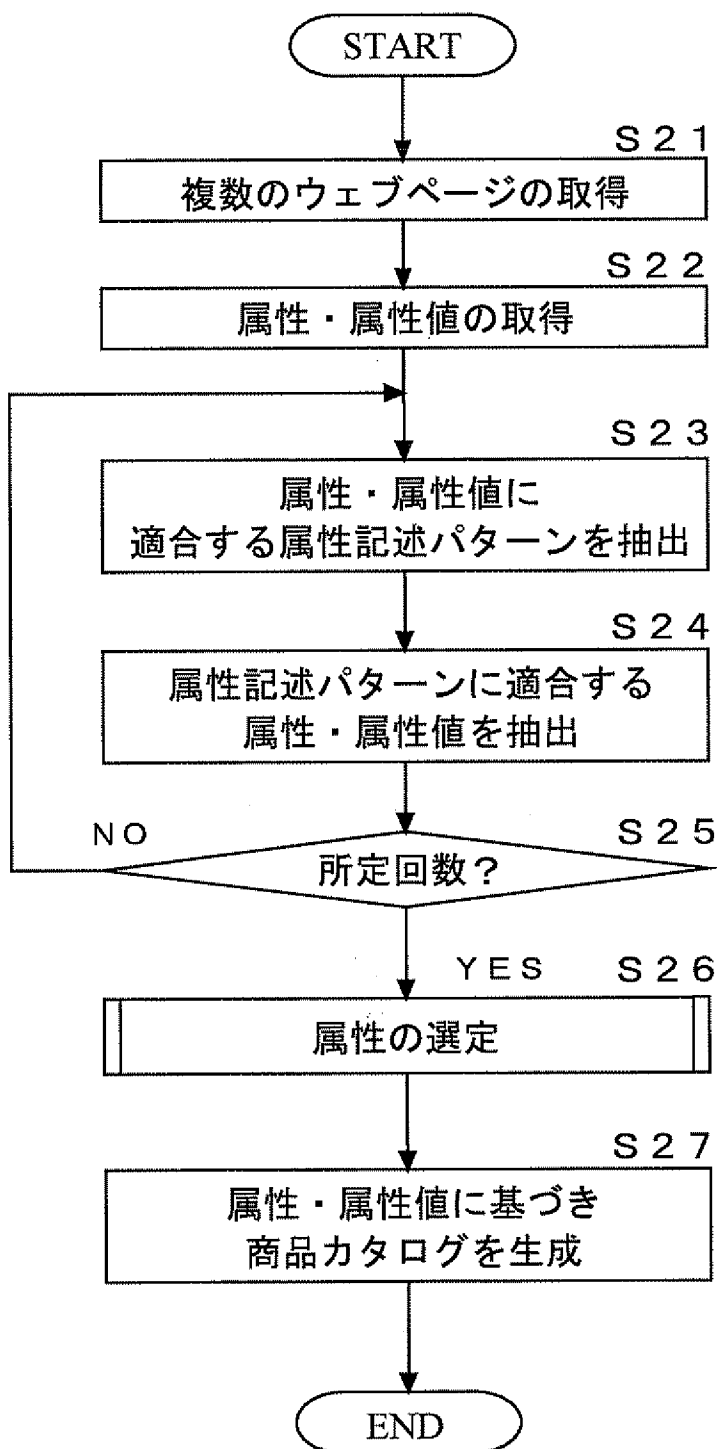
⋮
⋮
⋮

⋮
⋮
⋮

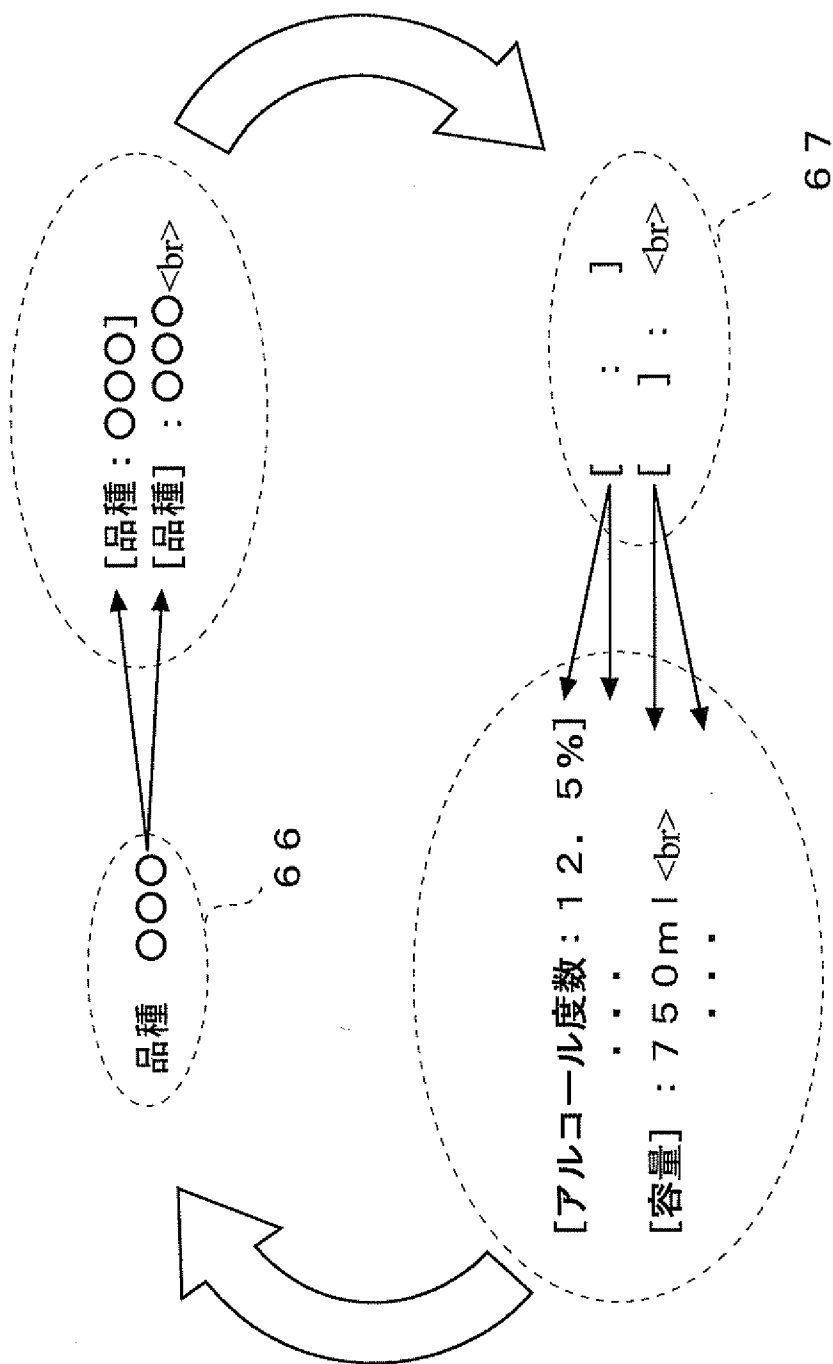
[図11]



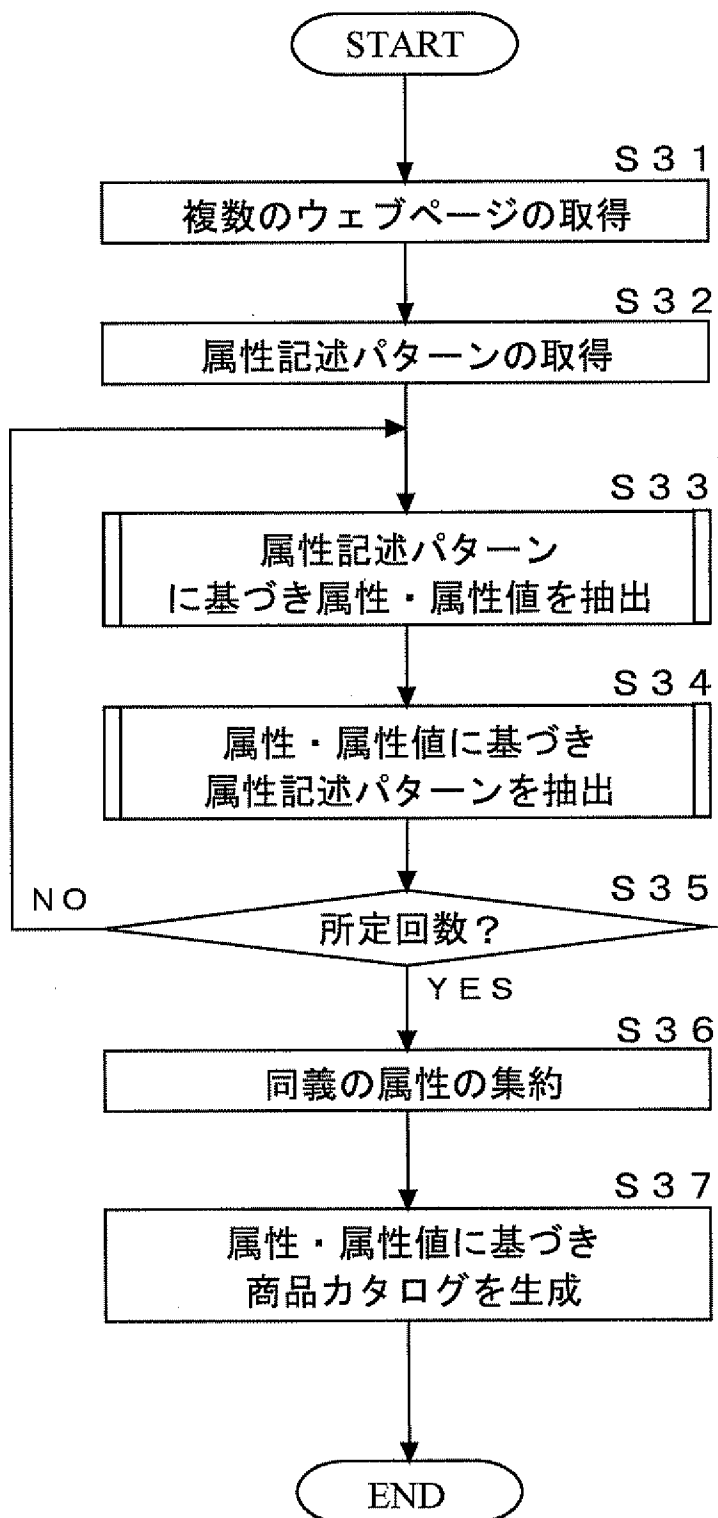
[図12]



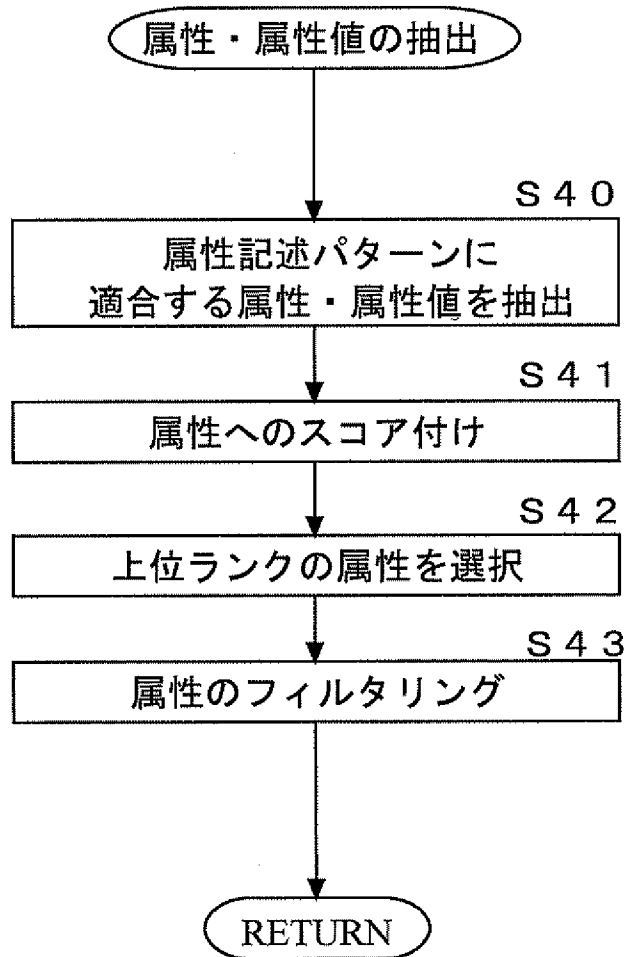
[図13]



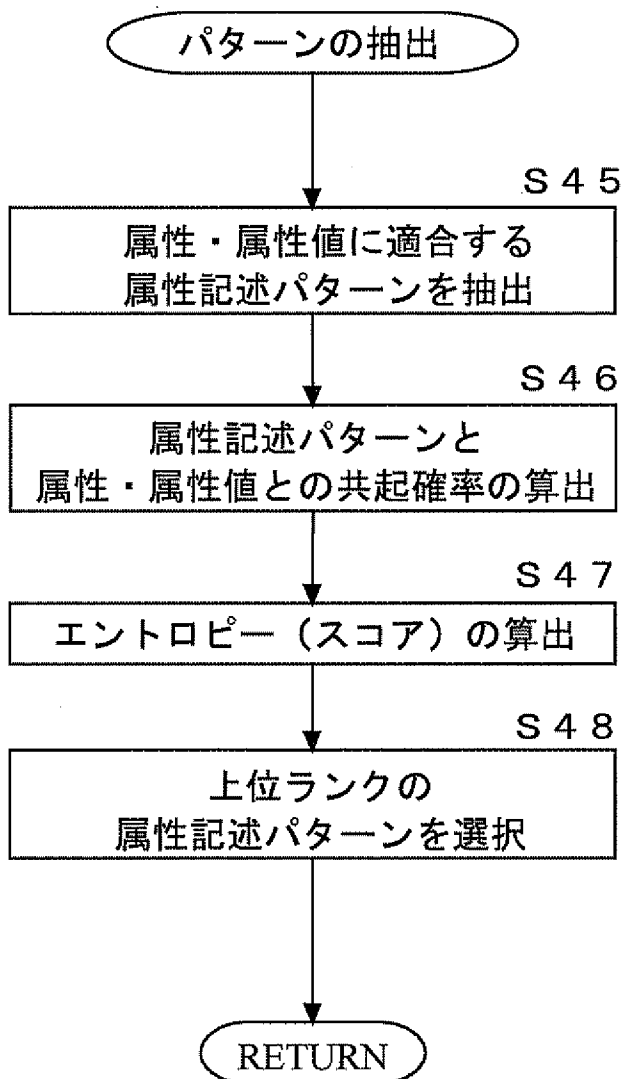
[図14]



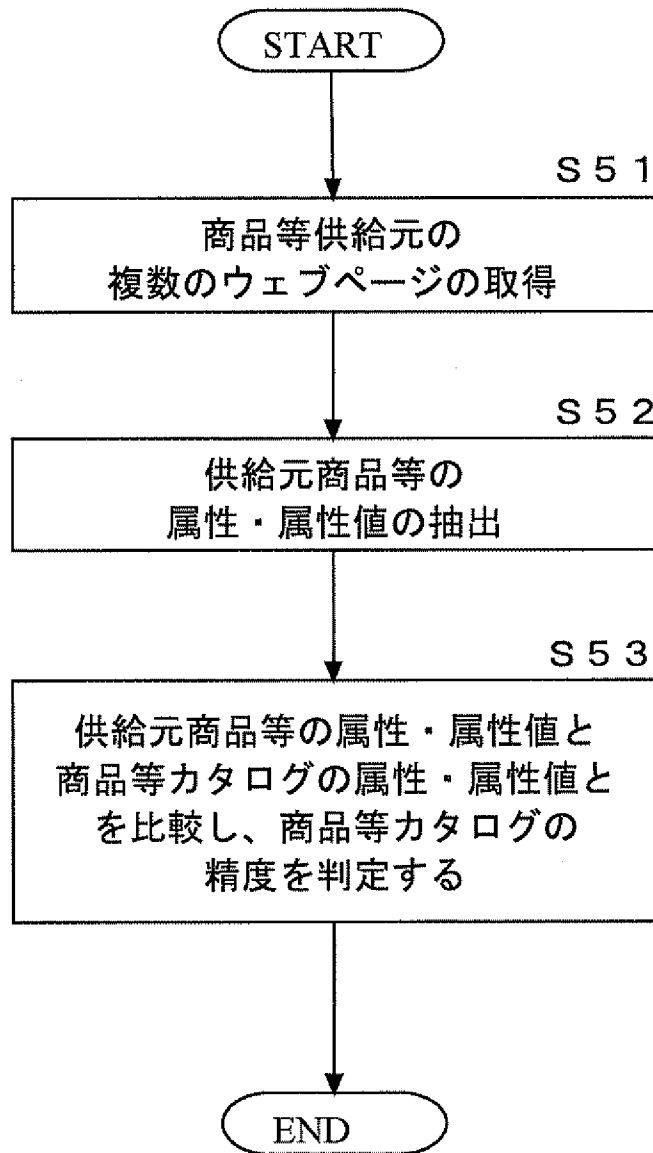
[図15]




[図16]



[図17]



[図18]

◆◆◆社	メルマガ	<input type="text"/>	検索
トップページ>ワイン>ドイツ			
ツエ〇・〇〇・△△		カタログ No.123456788	
		
		容量	750m l
		容器	びん
.....			
果実名	:	
収穫地	:	
醸造地	:	
.			
.			
.			
味のタイプ	:	

[図19]

製造年別カタログ

◆ 製造年：1995年

ワイン名	品種	生産者	価格	商品番号
●●ワイン	○○	△△	¥1,500	10461
▽▲ワイン	○●	▲△	¥2,900	12346

⋮

⋮

⋮

■▽ワイン	□○	○△	¥3,500	13354
-------	----	----	--------	-------

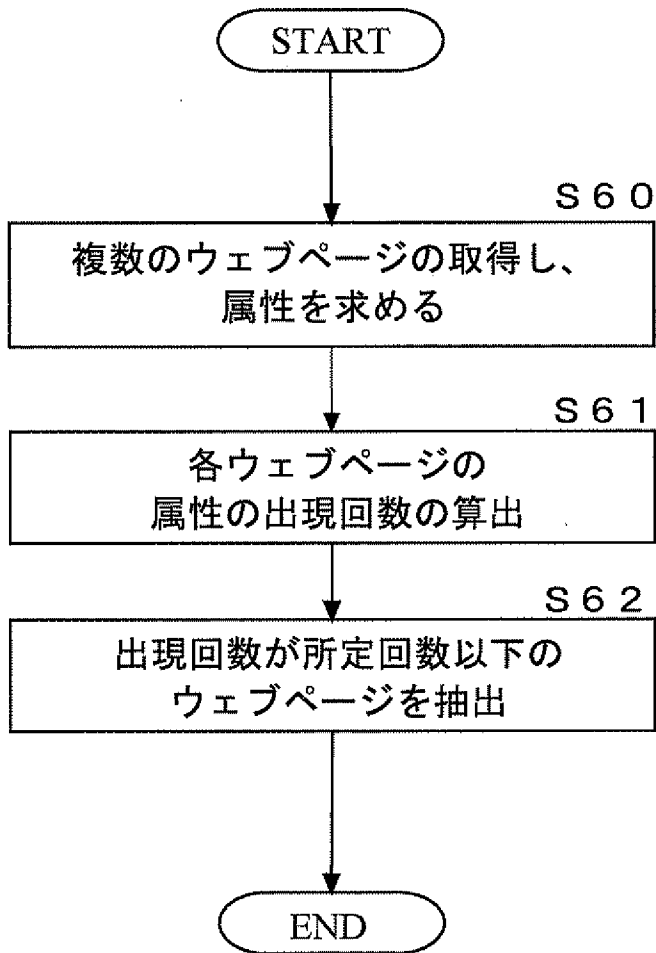
◆ 製造年：1996年

ワイン名	品種	生産者	価格	商品番号
●●ワイン	○○	△△	¥1,300	13478
▽▲ワイン	○●	▲△	¥2,500	15356

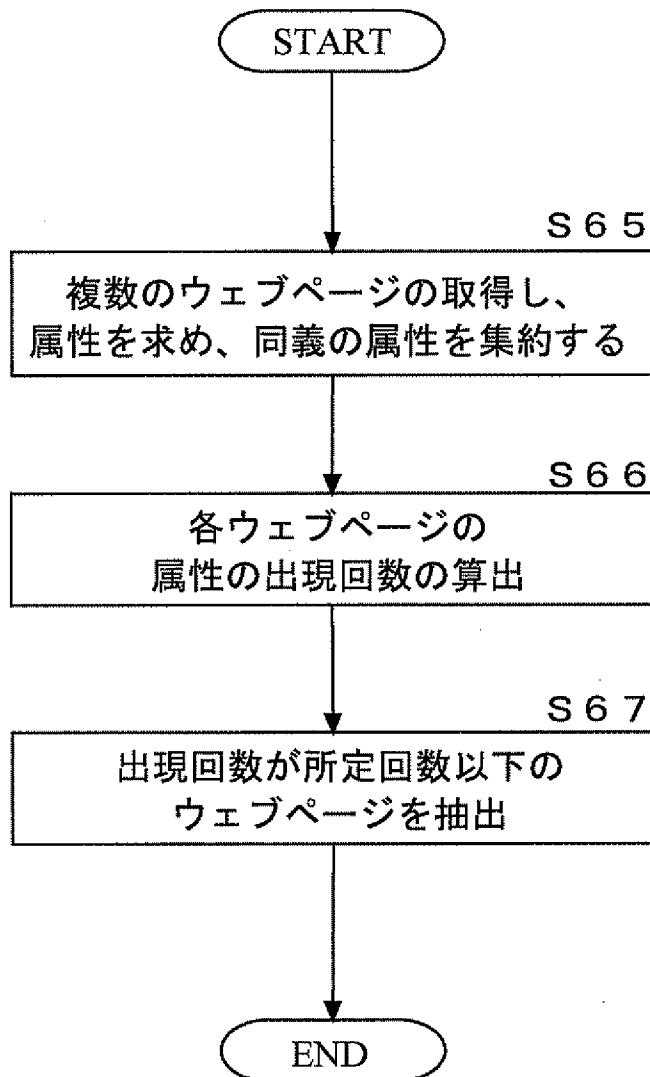
⋮

⋮

[図20]



[図21]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2011/054510

A. CLASSIFICATION OF SUBJECT MATTER

G06F17/30(2006.01) i, G06Q30/00(2006.01) i

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G06F17/30, G06Q30/00

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho	1922-1996	Jitsuyo Shinan Toroku Koho	1996-2011
Kokai Jitsuyo Shinan Koho	1971-2011	Toroku Jitsuyo Shinan Koho	1994-2011

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JP 2008-269106 A (Osaka Industrial Promotion Organization), 06 November 2008 (06.11.2008), entire text; all drawings (Family: none)	1-21
P, X	Satoshi SEKINE, "Shopping Site ni Okeru Shohin no Doitsusei, Ruijisei no Suitei Shuho", Proceedings of the 16th annual meeting of the Association for Natural Language Processing, 08 March 2010 (08.03.2010), pages 254 to 257	1-21

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&" document member of the same patent family

Date of the actual completion of the international search
15 March, 2011 (15.03.11)

Date of mailing of the international search report
29 March, 2011 (29.03.11)

Name and mailing address of the ISA/
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2011/054510

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P,X	Hiroki SAKACHI, "Shohin Page Karano Zokusei· Zokuseichi Chushutsu to Doitsu Shohin Clustering Shuho", Proceedings of the 16th annual meeting of the Association for Natural Language Processing, 08 March 2010 (08.03.2010), pages 371 to 374	1-21

A. 発明の属する分野の分類 (国際特許分類 (IPC))
 Int.Cl. G06F17/30(2006.01)i, G06Q30/00(2006.01)i

B. 調査を行った分野
 調査を行った最小限資料 (国際特許分類 (IPC))
 Int.Cl. G06F17/30, G06Q30/00

最小限資料以外の資料で調査を行った分野に含まれるもの
 日本国実用新案公報 1922-1996年
 日本国公開実用新案公報 1971-2011年
 日本国実用新案登録公報 1996-2011年
 日本国登録実用新案公報 1994-2011年

国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
X	JP 2008-269106 A (財団法人大阪産業振興機構) 2008.11.06, 全文, 全図 (ファミリーなし)	1-21
P, X	関根聡, ショッピングサイトにおける商品の同一性、類似性の推定 手法, 言語処理学会第16回年次大会発表論文集, 2010.03.08, p p. 254-257	1-21

C欄の続きにも文献が列挙されている。 パテントファミリーに関する別紙を参照。

<p>* 引用文献のカテゴリー 「A」特に関連のある文献ではなく、一般的技術水準を示すもの 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す) 「O」口頭による開示、使用、展示等に言及する文献 「P」国際出願日前で、かつ優先権の主張の基礎となる出願</p>	<p>の日の後に公表された文献 「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」同一パテントファミリー文献</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

国際調査を完了した日 15.03.2011	国際調査報告の発送日 29.03.2011
国際調査機関の名称及びあて先 日本国特許庁 (ISA/JP) 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官 (権限のある職員) 松田 直也 電話番号 03-3581-1101 内線 3599

C (続き) . 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
P, X	坂地泰紀, 商品ページからの属性・属性値抽出と同一商品クラスタリング手法, 言語処理学会第16回年次大会発表論文集, 2010.03.08, pp. 371-374	1-21