

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4637113号
(P4637113)

(45) 発行日 平成23年2月23日 (2011.2.23)

(24) 登録日 平成22年12月3日 (2010.12.3)

(51) Int. Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 340Z
 G06F 17/30 419A
 G06F 17/30 330B

請求項の数 8 (全 59 頁)

(21) 出願番号	特願2006-540102 (P2006-540102)	(73) 特許権者	000001007
(86) (22) 出願日	平成16年11月26日 (2004.11.26)		キヤノン株式会社
(65) 公表番号	特表2007-519086 (P2007-519086A)		東京都大田区下丸子3丁目30番2号
(43) 公表日	平成19年7月12日 (2007.7.12)	(74) 代理人	100076428
(86) 国際出願番号	PCT/AU2004/001676		弁理士 大塚 康德
(87) 国際公開番号	W02005/052810	(74) 代理人	100112508
(87) 国際公開日	平成17年6月9日 (2005.6.9)		弁理士 高柳 司郎
審査請求日	平成19年11月26日 (2007.11.26)	(74) 代理人	100115071
(31) 優先権主張番号	2003906611		弁理士 大塚 康弘
(32) 優先日	平成15年11月28日 (2003.11.28)	(74) 代理人	100116894
(33) 優先権主張国	オーストラリア (AU)		弁理士 木村 秀二
		(74) 代理人	100130409
			弁理士 下山 治
		(74) 代理人	100134175
			弁理士 永川 行光

最終頁に続く

(54) 【発明の名称】 階層データの好ましいビューを構築するための方法

(57) 【特許請求の範囲】

【請求項 1】

保持手段と、構築手段と、第一の生成手段と、入力手段と、第一の特定手段と、第二の特定手段と、第二の生成手段とを備える装置における方法であって、

前記保持手段が、複数のノードから構成される第一のツリーをデータベースに保持する保持工程と、

前記構築手段が、第一のコマンドに基づいて前記第一のツリーから第二のツリーを構築する構築工程と、

前記第一の生成手段が、前記第二のツリーに基づいて、(1) 前記第二のツリーの各ノードの生起する頻度を示す生起頻度テーブルと、(2) 前記第二のツリーの複数のノードの同時に生起する頻度を示す共起頻度テーブルとを生成する第一の生成工程と、

前記入力手段が、検索キーワードを入力する入力工程と、

前記第一の特定手段が、前記第一のツリーのノードの中から、前記入力工程で入力された前記検索キーワードに対応する複数のノードをヒットノードとして特定する第一の特定工程と、

前記第二の特定手段が、前記第一のツリーのノードの中から、前記第一の特定工程により特定された前記複数のヒットノードに共通する祖先ノードを特定する第二の特定工程と、

前記第二の生成手段が、前記生起頻度テーブルと、前記共起頻度テーブルと、前記祖先ノードとに基づいて、(a) 前記検索キーワードに対する検索結果として前記第一のツリ

10

20

一の一部を提示するための、(b)前記第二のツリーと共通のノードを含み、且つ、(c)前記複数のヒットノードに共通する祖先ノードをルートノードとする、コンテキストツリーを生成する第二の生成工程と、
を有することを特徴とする方法。

【請求項2】

前記第一のコマンドは複数存在し、

前記構築工程では、前記複数の第一のコマンドに基づいて、前記第一のツリーから複数の第二のツリーを構築することを特徴とする請求項1に記載の方法。

【請求項3】

前記共起頻度テーブルは、前記第二のツリーの複数のノードが同時に生起し、かつ子ノードを有しない葉ノードである頻度を示すテーブルを含むことを特徴とする請求項1又は2に記載の方法。

10

【請求項4】

前記共起頻度テーブルは、前記第二のツリーの複数のノードが同時に生起し、かつ子ノードが一つである頻度を示すテーブルを含むことを特徴とする請求項1又は2に記載の方法。

【請求項5】

前記装置は形成手段をさらに備え、

前記形成手段が、前記コンテキストツリーに基づいて、前記検索キーワードに対する検索結果として前記第一のツリーの一部を提示するために、第二のコマンドを形成する形成工程をさらに有することを特徴とする請求項1乃至4のいずれか1項に記載の方法。

20

【請求項6】

請求項1乃至5の何れか1項に記載の方法の各工程をコンピュータに実行させるためのプログラム。

【請求項7】

請求項6に記載のプログラムを格納したことを特徴とするコンピュータが読取可能な記憶媒体。

【請求項8】

複数のノードから構成される第一のツリーをデータベースに保持する保持手段と、

第一のコマンドに基づいて前記第一のツリーから第二のツリーを構築する構築手段と、

前記第二のツリーに基づいて、(1)前記第二のツリーの各ノードの生起する頻度を示す生起頻度テーブルと、(2)前記第二のツリーの複数のノードの同時に生起する頻度を示す共起頻度テーブルと、を生成する第一の生成手段と、

30

検索キーワードを入力する入力手段と、

前記第一のツリーのノードの中から、前記入力手段により入力された前記検索キーワードに対応する複数のノードをヒットノードとして特定する第一の特定手段と、

前記第一のツリーのノードの中から、前記第一の特定手段により特定された前記複数のヒットノードに共通する祖先ノードを特定する第二の特定手段と、

前記生起頻度テーブルと、前記共起頻度テーブルと、前記祖先ノードとに基づいて、(a)前記検索キーワードに対する検索結果として前記第一のツリーの一部を提示するための、(b)前記第二のツリーと共通のノードを含み、且つ、(c)前記複数のヒットノードに共通する祖先ノードをルートノードとする、コンテキストツリーを生成する第二の生成手段と、

40

を備えることを特徴とする装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報検索(retrieval)の一般分野に関し、特に、関連性の高いデータを大規模な階層データの情報源から自動的に識別および検索することに関する。

【背景技術】

50

【0002】

XML (Extensible Markup Language、拡張可能なマーク付け言語)は、情報を保存および交換するための一般的な階層フォーマットにますますなりつつある。XMLの階層的な性質により、XMLは、データ・オブジェクト間の関係を捉えるための優れた手段となる一方で、キーワード検索をより困難にもしている。

【0003】

キーワード検索は、XMLなどの構造化されたデータ・フォーマットを扱うときには特に重要であるが、その理由は、これによってユーザが特定のキーワードを、データの内部構造を知ることを必要とせず迅速に見つけられるようになることにある。これはXMLで作業を行うときに難問となるが、それは、キーワード検索の結果を提示するための最適のまたは明らかに好ましい方法がないためである。伝統的な構造化されていないテキストの環境では、データ・システムから、通常、ユーザに対して、見つかったキーワードをその付近にある他のテキストと一緒に提示される。「ヒット」が1つより多い場合、隣接するテキストにより、ヒット同士を区別するための有用なコンテキストを提供され、それにより、ユーザは、関連性の最も高いヒットの選択をユーザの必要に応じて行えるようになる。

10

【0004】

構造化されているXMLの環境では、互いに関係するデータはXML文書内部のいくつかのばらばらな(disjoint)場所にある可能性があるため、「隣接するデータ」(neighbouring data)のはっきりした概念がない。このため、キーワード検索においてヒットに適したコンテキストを識別または構築することが難しい。したがって、既存のたいていのXMLベースのシステムでは、ある文書の内部でキーワード・ヒットが生じた場合、(XML文書からなるある集まりの中から)単にそのXML文書全体を返し、実質的にその文書全体にそのヒットに対するコンテキストの役割をさせている。これは、文書が大きく、ユーザがその内容のすべてを見ることに興味のないときには、望ましくない。

20

【0005】

現実的なデータソース、特にデータベースは、しばしば、ユーザがとにかく一度で見たいと思うよりもずっと多くのデータを含んでいる。たとえば、メール注文の店にあるデータベースには、その製品、顧客、納入業者、運送業者、ならびに過去および未発送の注文のすべてについての詳細が含まれているかもしれない。店員は、特定の製品の現在の在庫量を見たいと希望することもあるし、ある顧客に対する注文のステータスを確認したいと思うこともある。他方、店のマネージャーが、何ヶ月にもわたって特定の製品の総売り上げの変動を見たいと希望することがあるかもしれない。こうした場合のそれぞれで、さらに別の関連性のないデータが雪崩のように提示されるとすれば、それはあまりにもユーザを惑わせることになる。さらに、ユーザは、そのデータベースの構造に詳しくない限り、通常は、自分の興味がある情報を特定することができないはずである。

30

【0006】

関連性の高いデータだけを提供するための伝統的な方法は、システム管理者など、データソースの構造に詳しい誰かの準備した、あらかじめ作成済みの「ビュー」の使用によるものである。各ビューは、そのデータソースの何らかのサブセットを引き出してまとめるものであり、はっきりした目的に合わせて作られている。先に挙げた例では、店員は「在庫量」ビューまたは「注文ステータス」ビューを調べるであろうし、マネージャーは「売上」ビューを持ち出すことになる。

40

【0007】

このあらかじめ作成済みのビューを使用するというアプローチは、ありそうな使用法の想定状況をすべて見越しておけるときには十分であるが、キーワード検索には適していない。キーワード検索の操作では、ユーザは、1つまたは複数のキーワードを入力し、システムは、(ANDのプール値キーワード検索操作を仮定すると)そのキーワードすべての出現を含むデータ・セットまたはビューによって応答する。XMLなどの階層的な環境で

50

は、キーワード・ヒットは、階層内の異なる場所にあるいくつかのデータ項目の中に生じる可能性がある。ユーザの与える可能性のある可能なキーワードの組合せすべてを見越しておくことは実現可能でない以上、階層内のどこでヒットが生じるかをあらかじめ判定することは可能でない。したがって、検索の想定状況すべてをまかなうようなあらかじめ作成済みのビューを提供することは可能ではない。

【0008】

これに類したキーワード検索の問題点は、関係データベースの環境でも存在する。関係データベースは、そのプライマリおよびフォーリンのキーによって併合されたテーブルを含み、ここで各テーブルは、あるエンティティに対する属性の値からなる n -組 (n -tuple) をそれぞれが表す複数の行を含む。Hristidis, V. and Papakonstantinou, Y., 「DISCOVER: Keyword Search in Relational Databases」、Proceedings of the 28th VLDB Conference、2002年で記述されている、関係データベースにおけるキーワード検索に対する伝統的な解決策は、キーワード・ヒットすべてを含む併合済みテーブル全体にわたって併合される行の最小のネットワークである、最小併合ネットワーク (minimal joining network) を返すことである。このアプローチの問題点は、キーワード・ヒットがデータベース・テーブルのある行の中のどこかで生じた場合、その行全体がそのヒットに対するコンテキストとして返されるという点で、行が最小のデータ「チャンク」として効果的に扱われることである。このことから、極端な量のデータがユーザに提示されることになる可能性があるが、これは、典型的な関係データベースのテーブルが、しばしば、ユーザには普通興味のない多くの列を含んでいるためである。

【0009】

さらに、上の技法を、XMLなどの階層的なデータ構造に適合させると、コンテキストの情報が不十分になる可能性がある。階層的な環境では、関係するデータは、階層内の異なるレベルで格納されている可能性があり、このためしばしば、親または祖先のノードあるいはその子からは、キーワード検索に対する非常に有用なコンテキストが、たとえ最小データ・セットには含まれないかもしれない場合でも、提供される可能性がある。

【0010】

階層データにおけるキーワード検索の問題に取り組むために、いくつかの試みがなされてきている。Florescu, D. 「Integrating Keyword Search into XML Query Processing」、Ninth International World Wide Web Conference、2000年5月では、構造化問い合わせ言語をキーワード検索演算子 `contains` で拡張する方法が開示されている。この演算子では、指定したサブツリーがいくつかの指定したキーワードを含む場合に、評価が真となる。ユーザは、この演算子を、問い合わせを構築するときに使用して、望ましくないデータをふるい除けることができる。この有用な機能では、ユーザは、所与のサブツリー内部のヒットしたキーワードの正確な場所を指定する必要はないが、返されるデータの正確なフォーマットをやはり検索問い合わせの中で指定する必要があり、このため、データソースの構造に詳しいことが必要であるため、十分ではない。言い換えると、フリー・テキストのキーワード検索は、ユーザがデータソース全体を検索の結果として受け入れる気がない限り、依然として可能ではないのである。

【0011】

XMLデータソースにおけるキーワード検索に対する既存の別のアプローチでは、ユーザが、スキーマ要素からなる所与のリストから、返されるデータのルート・ノードに相当する要素を選択する必要がある。キーワード・ヒットが生じるのが、選択したスキーマ要素の表すデータ要素の子孫ノード内である場合、ヒットしたキーワードを含む、そのデータ要素の下のサブツリー全体がユーザへと返される。このアプローチは、ユーザの介入を必要とするため、扱いにくい。さらに、ユーザは、サブツリー全体を、たとえそのユーザに興味のないデータを含むかもしれない場合でも、受け入れるよう強いられる。

【0012】

したがって、階層的なデータの環境における1組の関連性の高いデータの判定を、ユーザの介入またはユーザの事前の階層データの構造の知識を必要としない、キーワードの任意の組合せを含むキーワード検索操作に回答して行うための方法が必要とされている。

【非特許文献1】Hristidis, V. and Papakonstantinou, Y., 「DISCOVER: Keyword Search in Relational Databases」、Proceedings of the 28th VLDB Conference、2002年

【非特許文献2】Florescu, D. 「Integrating Keyword Search into XML Query Processing」、Ninth International World Wide Web Conference、2000年5月

【特許文献1】オーストラリア特許出願第2003204824号

【特許文献2】米国特許出願第10/465222号

【発明の開示】

【発明が解決しようとする課題】

【0013】

既存の諸方法の1つまたは複数の不利な点を克服または少なくとも改善することが、本発明の一目的である。

【課題を解決するための手段】

【0014】

上記の目的を達成する本発明に係る方法は、

保持手段と、構築手段と、第一の生成手段と、入力手段と、第一の特定手段と、第二の特定手段と、第二の生成手段とを備える装置における方法であって、

前記保持手段が、複数のノードから構成される第一のツリーをデータベースに保持する保持工程と、

前記構築手段が、第一のコマンドに基づいて前記第一のツリーから第二のツリーを構築する構築工程と、

前記第一の生成手段が、前記第二のツリーに基づいて、(1)前記第二のツリーの各ノードの生起する頻度を示す生起頻度テーブルと、(2)前記第二のツリーの複数のノードの同時に生起する頻度を示す共起頻度テーブルとを生成する第一の生成工程と、

前記入力手段が、検索キーワードを入力する入力工程と、

前記第一の特定手段が、前記第一のツリーのノードの中から、前記入力工程で入力された前記検索キーワードに対応する複数のノードをヒットノードとして特定する第一の特定工程と、

前記第二の特定手段が、前記第一のツリーのノードの中から、前記第一の特定工程により特定された前記複数のヒットノードに共通する祖先ノードを特定する第二の特定工程と

前記第二の生成手段が、前記生起頻度テーブルと、前記共起頻度テーブルと、前記祖先ノードとに基づいて、(a)前記検索キーワードに対する検索結果として前記第一のツリーの一部を提示するための、(b)前記第二のツリーと共通のノードを含み、且つ、(c)前記複数のヒットノードに共通する祖先ノードをルートノードとする、コンテキストツリーを生成する第二の生成工程と、

を有することを特徴とする。

【0021】

ここで、本発明の少なくとも1つの実施形態を、図面を参照して説明することにする。

【発明を実施するための最良の形態】

【0022】

以下に、図面を参照して、この発明の好適な実施の形態を例示的に詳しく説明する。ただし、この実施の形態に記載されている構成要素はあくまで例示であり、この発明の範囲

10

20

30

40

50

をそれらのみ限定する趣旨のものではない。

【0023】

本開示により、階層的なデータの環境における1組の関連性の高いデータの判定を、1つまたは複数のキーワードを含むキーワード検索操作に回答して行うための方法が提供される。好ましい一実施形態は、階層的なデータ構造におけるデータの好ましいビューの構築を、データが過去のエピソードでどのようにアクセスされたかに基づいて行うベイズ確率に基づく方法を含む。より具体的には、この方法では、データ項目の対とベクトルの間の過去の同時生起の頻度を利用して、あるデータ項目がいくつかの他の義務的データ項目にとって関連性が高いことの確率を計算する。通常、義務的データ項目は、キーワード・ヒットを含むものであり、このため、ユーザに対してキーワード検索結果に入れて返さなければならぬ。義務的でないデータ項目は、義務的データ項目にとって関連性が高いことの確率が高い場合には、検索結果に入れて返されてキーワード・ヒットに対するコンテキストとして役立つことが見込まれる。

10

【0024】

ここに開示する配置の、伝統的なあらかじめ作成されたビューに基づくアプローチに対する顕著な特徴は、前者が、単に既存の保存されたビューを返すのではなく、新しいビューを合成できることである。このため、そのような配置は、任意のキーワード組合せを含むキーワード検索操作を扱うことができ、またビューは、動的に生成されるため、あらかじめ作成されたビューの固定のプールから得られるものよりも、個々の操作によりよく合わせることができる。

20

【0025】

ここに開示する方法では、通常、いくつもの代替のビューを構築し、そのビューがユーザにとってどのくらい興味があるものかを示すスコアを、各ビューに対して割り当てる。ある実装形態では、構築したもののうちでスコアの最も高い単一のビューがユーザに返される。ある代替実装形態では、最高から最低へとそのスコアに従ってソートされたビューのリストが返される。

【0026】

ここに開示する方法は、キーワード検索がその第一の動機ではあるが、ともに「Methods for Interactively Defining Transformations and for Generating Queries by Manipulating Existing Query Data」という名称の2003年6月19日に出願されたオーストラリア特許出願第2003204824号、およびこれに対応する、2003年6月20日に出願された米国特許出願第10/465222号に記載のものなどの、階層データの提示の方法を向上させるのに使用することもできる。その公開(publication)では、最も適切な提示タイプ(テーブル、グラフ、プロット、ツリーなど)の選択を階層的なデータソースの構造および内容に基づいて行うための方法が開示されている。この方法は、提示タイプの選択の前に、本開示の好ましい一実施形態を、表示に最も好ましいデータソースのサブセットを自動的に選択する手段として組み込むことによって向上させることができる。データの好ましいサブセットだけをこのようにして表示することがしばしば有用であるが、これは、階層的なデータソースが、しばしば、ユーザに普通興味のあるはずのものよりも多くの情報を含むためであり、このため、本発明の好ましい実施形態などの「興味のない」データをふるい除ける方法は、ユーザの経験をより満足の行く生産的なものにするのに役立つ可能性がある。

30

40

【0027】

以下にある説明のいくつかの部分は、明示的にまたは暗黙に、コンピュータ・メモリ内部のデータに対する操作のアルゴリズムおよび記号表現によって提示してある。こうしたアルゴリズムの説明および表現は、データ処理の当業者がその作業の実質を他の当業者に最も効率よく伝えるのに使用する手段である。アルゴリズムは、ここでは、また一般的には、所望の結果をもたらす自己整合的な一連のステップと考えられている。ステップは、物理量の物理的操作を必要とするものである。こうした量は、必ずではないが、普通、保

50

存、転送、組合せ、比較、および他の方法で操作するのに堪える電気または磁気の信号の形式をとる。時折、主に一般的に使用されているという理由で、こうした信号をビット、値、要素、記号、キャラクタ、項、数などと呼ぶと好都合であることがわかっている。

【0028】

しかし、上のおよび類似の用語は、適切な物理量に関連づけられるべきであり、そうした量に適用される便宜的なラベルにすぎないことを銘記すべきである。そうでないと特に述べない限り、また以下から明らかなように、本明細書を通して、「スキャンする」、「計算する」、「判定する」、「置き換える」、「生成する」、「初期化する」、「出力する」などの用語を使用する議論が指すのは、コンピュータ・システムのレジスタおよびメモリ内部の物理（電子的）量として表現されたデータを、コンピュータ・システムのメモリまたはレジスタあるいはそのような他の情報の保存、伝達、または表示の装置内部の物理量として同様に表現される他のデータへと操作または変換するコンピュータ・システムまたは類似の電子的装置の動作およびプロセスであることが理解されよう。

10

【0029】

また、本明細書では、この方法の操作を行うための機器の開示を行う。そのような機器は、必要な目的向けに特別に構築することができ、またはコンピュータ内に格納したコンピュータ・プログラムによって選択的に活性化または再構成される汎用コンピュータまたは他の装置を含んでもよい。本明細書で提示するアルゴリズムおよび表示は、特定のどのようなコンピュータまたは他の機器とも固有の関係をもたない。様々な汎用の機械を、本明細書中の教示によるプログラムとともに使用することができる。あるいは、必要な方法ステップを行うためにより特化した機器を構築することが適切であることもある。従来の汎用コンピュータの構造は、下の説明の中で行うことにする。

20

【0030】

さらにまた、本明細書では、この方法の操作を行うためのコンピュータ・プログラムを含むコンピュータ可読媒体の開示を行う。コンピュータ可読媒体は、本明細書では、送信元とあて先の間でコンピュータ・プログラムを通信で送るためのどのような伝送媒体も含むと解釈する。伝送媒体は、磁気もしくは光ディスク、メモリ・チップなどのストレージ装置、または汎用コンピュータとのインターフェースを行うのに適した他のストレージ装置を含むことができる。また、伝送媒体は、インターネット・システムの形で例示されるような配線接続（hard-wired）の媒体、またはGSM移動体電話システムの形で例示されるような無線媒体を含んでもよい。コンピュータ・プログラムは、特定のどのようなプログラミング言語およびその実装形態にも限定されるものではない。様々なプログラミング言語およびそのコーディングは、本明細書に含まれる開示の教示を実装するのに使用できることが理解されよう。

30

【0031】

添付の図面の1つまたは複数のいずれかにおいて、符号（reference numerals）の同じステップおよび/または特徴への参照を行っている場合、これらのステップおよび/または特徴は、この説明の目的にとっては、それに反する意図を示していない限り、機能または操作は同じである。

【0032】

この一般にはキーワード検索、特に階層的なデータ構造の構築の諸方法の実施は、好ましくは、図1から36のプロセスを、コンピュータ・システム3700内部で実行されるアプリケーション・プログラムなどのソフトウェアとして実装できる、図37に示すものなどの汎用コンピュータ・システム3700を用いて行う。特に、キーワード検索の諸ステップは、コンピュータの実行するソフトウェアの中の命令によってもたらされる。この命令は、それぞれが1つまたは複数の特定のタスクを行うための、1つまたは複数のコード・モジュールとして形成することができる。また、このソフトウェアは、第1の部分が検索の方法を行い、第2の部分が第1の部分とユーザとの間のユーザ・インターフェースを管理する、2つの別々の部分に分割することができる。次いで、このソフトウェアは、たとえば、下で説明するストレージ装置を含むコンピュータ可読媒体に格納することがで

40

50

きる。このソフトウェアは、そのコンピュータ可読媒体からコンピュータへとロードされ、次いで、そのコンピュータによって実行される。そのようなソフトウェアまたはコンピュータ・プログラムがその上に記録してあるコンピュータ可読媒体が、コンピュータ・プログラム製品である。そのコンピュータ内のコンピュータ・プログラム製品の使用により、好ましくは、キーワード検索および階層的なデータ構造の構築のための有利な機器がもたらされる。

【0033】

コンピュータ・システム3700は、コンピュータ・モジュール3701、キーボード3702やマウス3703などの入力装置、プリンタ3715、ディスプレイ装置3714、およびラウドスピーカ3717を含む出力装置によって形成される。モデム(Modem、Modulator-Demodulator)トランシーバ装置3716は、コンピュータ・モジュール3701で使用して、たとえば、電話回線3721または他の機能媒体(functional medium)を介して接続可能な、通信ネットワーク3720との間の通信が行われる。モデム3716は、インターネット、および、LAN(Local Area Network、構内通信網)やWAN(Wide Area Network、広域通信網)など、他のネットワーク・システムへのアクセスを得るのに使用することができ、一部の実装形態では、コンピュータ・モジュール3701へと組み込むこともできる。

【0034】

コンピュータ・モジュール3701は、通常、少なくとも1つのプロセッサ・ユニット3705と、たとえば半導体RAM(random access memory)およびROM(read only memory)から形成される、メモリ・ユニット3706とを含む。また、モジュール3701は、ビデオ・ディスプレイ3714およびラウドスピーカ3717へと結合されるオーディオ・ビデオ・インターフェース3707、キーボード3702およびマウス3703、および自由選択でジョイスティック(図示せず)に対する入出力インターフェース3713、ならびにモデム3716およびプリンタ3715に対するインターフェース3708を含むいくつかの入出力(I/O)インターフェースを含む。一部の実施形態では、モデム3716を、コンピュータ・モジュール3701内部、たとえばインターフェース3708内部に組み込んでよい。ストレージ装置3709が提供されており、通常、ハード・ディスク・ドライブ3710およびフロッピー・ディスク・ドライブ3711を含む。また、磁気テープ・ドライブ(図示せず)を使用してもよい。CD-ROMドライブ3712は、通常、データの非揮発性の情報源として提供されている。コンピュータ・モジュール3701のコンポーネント3705から3713は、通常、通信を、相互接続されたバス3704を介して、また当業者に知られているコンピュータ・システム3700の従来の利用モード(mode of operation)になるような仕方で行う。ここで説明する配置を実施できるコンピュータの例としては、IBM-PCおよびその互換機、Sun Sparcstations、またはそれから発展した類似のコンピュータ・システムがある。

【0035】

通常、アプリケーション・プログラムは、ハード・ディスク・ドライブ3710上にあり、プロセッサ3705により、その実行中に読み込みおよび制御が行われる。このプログラムおよびネットワーク3720からフェッチされるどのようなデータの間接ストレージも、半導体メモリ3706を用いて、おそらくハード・ディスク・ドライブ3710と協力して達成される。一部の場合には、このアプリケーション・プログラムは、CD-ROMまたはフロッピー・ディスク上にエンコードされてユーザへと供給され、対応するドライブ3712または3711を介して読み込まれるが、または代替方法として、ユーザがネットワーク3720からモデム装置3716を介して読み込むのでもよい。またさらに、このソフトウェアのコンピュータ・システム3700へのロードは、他のコンピュータ可読媒体から行うことも可能である。本明細書で使用する「コンピュータ可読媒体」という用語は、命令および/またはデータを実行および/または処理のためにコンピュータ

10

20

30

40

50

・システム 3700 へと提供するにあずかるどのようなストレージ媒体または伝送媒体も指す。ストレージ媒体の例としては、フロッピー・ディスク、磁気テープ、CD-ROM、ハード・ディスク・ドライブ、ROM または集積回路、光磁気 (magneto-optical) ディスク、または PCMCIA カードなどのコンピュータ可読カードなどがあり、そのような装置がコンピュータ・モジュール 3701 の内部または外部にあることを問わない。伝送媒体の例としては、無線または赤外線伝送チャネルならびに別のコンピュータまたはネットワーク化装置へのネットワーク接続、電子メール転送および Web サイト上に記録された情報を含むインターネットまたはイントラネットなどがある。

【0036】

階層的な環境におけるキーワード検索は、その 1 つまたは複数のキーワードが生起する階層的なデータ構造内のノードまたは要素を識別し、次いで、他のどのような要素がそのキーワードにとって関連性が高いかを判定することを含む。典型的なキーワード検索の想定状況では、ユーザへと提示される結果のデータは、第 1 のデータ構造から抽出され、検索キーワードのすべてまたは一部およびこうしたキーワードにとって関連性が高いと考えられる他のデータを含む第 2 の階層的なデータ構造である。そのような、キーワード検索操作の結果としてユーザへと提示される階層的なデータ構造を、コンテキストツリー (context tree) と呼ぶ。

【0037】

XML ではしばしばあることであるが、検索されている階層的なデータに支配的な (governing) スキーマがあるとき、一般に有利なのは、そのスキーマのレベルで動作する関連性の高いデータを識別するための方法を使用することである。すなわち、スキーマ表現内部の要素が、解析されて、検索キーワードにとって関連性が高いかが判定される。次いで、関連性の高いスキーマ要素によって総体的に表現されるデータソース内のデータ項目のすべてのインスタンスが、キーワード検索操作の結果としてユーザへと返される。XML では、支配的なスキーマは、それ自体が別の階層的なデータ構造である、XML Schema の形をとることが可能である。XML Schema では、関連する XML データの構造、その XML データ内の要素および属性のリスト、およびそれらの親子関係が指定される。XML Schema 内の要素または属性のそれぞれは、通常、XML データ内の要素および属性の多くのインスタンスを表現しているため、XML Schema は、潜在的にずっと小さなデータ構造であり、このため、より効率よく解析することができる。

【0038】

しばしば、キーワードの検索を、1 つより多くの階層的なデータソースの中で行うことが望ましい。階層的なデータソースそれぞれは、それ自体はツリー構造をしているが、複数のデータソースを考慮すると、結果となるデータ構造はより一般的な形式をとる可能性がある。データベースの環境で不可避免的に生じるそのようなある形式を図 1 に示している。この構造は、本質的に、共有ノードのあるいくつものツリーを含み、ここで、ツリーそれぞれは、別々の階層的なデータソースのスキーマを表現しており、共有ノードは、内容が複数のデータソースにわたるデータ・ビューの結果である。具体的には、図 1 の点線の矩形 1005 および 1010 は、第 1 および第 2 のデータソースのスキーマをそれぞれ表し、ノード 1015 は、第 1 のデータソースからのノード 1020 および 1025 と第 2 のデータソースからのノード 1030 を一緒にするデータ・ビューのルート・ノードである。図 1 の複数の共有ツリーの構造は、本明細書ではこれをスキーマ・グラフ (schema graph) と呼ぶが、有向非巡回グラフの特別な形式であり、どの 2 つのノード間にも多くとも 1 つの有向パスがあるという重要な特徴をもつ。たとえば、ノード 1015 からノード 1035 へは有向パスが 1 つしかなく、このパスはノード 1020 を通過する。

【0039】

スキーマ・グラフは、好ましくは、キーワード検索操作の前に構築され、階層的なデータソースの初めは互いに素な個々のツリー構造をしたスキーマからなる。次いで、こうし

10

20

30

40

50

たスキーマツリーは、1つより多くのデータソースにわたるデータ・ビューが作成されると、併合される。データ・ビューは、通常、(XMLの環境におけるXQueryなどの)問い合わせを含み、データベース管理者またはユーザが作成することができる。どちらの場合でも、データベース・システムは、好ましくは、こうした問い合わせをログまたは記録をそのストレージ装置で行う。スキーマ・グラフの構築中、ログ済みの問い合わせそれぞれのスキーマ表現が作成され、そのスキーマ・グラフへと挿入される。こうすると、2つ以上の別々のスキーマツリーの併合が、挿入されるスキーマが図1に示すこれらのツリーからのノードを含む場合には、行われることになる。スキーマ・グラフへと挿入されるスキーマが、ただ1つのデータソースからのノードしか含まないこともあり得るのであり、その場合には、別々のスキーマツリーの併合は生じない。その代わりに、挿入操作によって、単に、新しいノードがスキーマ・グラフ内の単一のスキーマツリーからの既存のノードへと追加およびリンクされることになる。

10

【0040】

スキーマ・グラフの更新は、絶えず、新しい問い合わせのログのたびに行うこともでき、または新しい問い合わせおよびデータ・ビューがある期間にわたって集められた後、1つまたは複数の機会に行うこともできる。スキーマ・グラフの更新の頻度にかかわらず、キーワード検索操作が開始されると、その操作の時点で現在のものである(curren t)スキーマ・グラフは、ユーザへと返されるデータ・ビューを判定するのに使用されるものである。この文書の残りの部分では、「スキーマ・グラフ」という用語は、キーワード検索操作が行われる時点で現在のものである(curren t)スキーマ・グラフを指す。

20

【0041】

単一のデータソースの場合でのように、複数の階層的なデータソースの内部でのキーワード検索は、まず、「ヒット」ノードと呼ばれる、検索キーワードの見つかったスキーマ・グラフの内部のノードを識別し、次いで、「コンテキスト」ノードと呼ばれる、そのヒット・ノードにとって関連性の高いノードを識別することを含む。次いで、ヒット・ノードおよびコンテキスト・ノードを含むデータ構造が構築され、ユーザへと提示される。ヒット・ノードは1つより多くのデータソースにあり得るため、ユーザへと提示される結果のデータ構造は、複数のデータソースにわたる可能性がある。また、結果のデータ構造は、その意図している適用例は階層的なデータの環境においてであるため、好ましくは、ツリー構造をしている。

30

【0042】

図4に、キーワード検索操作の好ましい構成4000および一般化された利用モード(mode of operation)を示している。構成4000は、ネットワーク内で一緒に接続された、PCクライアント4005、データ・サーバ4010、データベース4015、キーワード検索クライアント4025、およびインデックス・サーバ4030を含む。

【0043】

装置4005、4010、4025、および4030のそれぞれは、通常、システム3700などの対応する汎用コンピュータ・システムによって形成され、それぞれはネットワーク3720によってリンクされるが、これを図4にもっぱら概念的に示している。この概念的な図示を使用して、様々な装置4005、4010、4025、および4030の間の、ネットワーク3720の両側で生じるデータ・フローの煩雑でない表現を提供している。様々な装置4005、4010、4025、および4030は、必要、適当、または好都合なときには、数のより少ない別個のコンピュータ・システム3700へとまとめることができる。たとえば、一部の実装形態では、サーバ4010と4030を組み合わせることで1つのコンピュータ・システム3700にし、クライアント4005と4025を組み合わせることで別のコンピュータ・システム3700にし、それらのシステム3700はネットワーク3720でリンクされるようにすることが好都合である可能性がある。

40

【0044】

50

データベース4015に格納されているデータへのアクセスは、通常、PCクライアント4005でブラウズを行うユーザによって行われる。クライアント4005内で動作するブラウズ用アプリケーションからは、コマンドの発行が、好ましくは、XQuery4006の形式で行われ、次いで、これがデータ・サーバ4010へと送信される。XQuery4006は、それぞれ、ログ4020に記録され、データ・サーバ4010によって解析され、その後、要求されたデータ4007がデータベース4015からフェッチされ、PCクライアントのユーザ4005へと送付される。ある時点、好ましくは十分な量のXQuery4006のログが行われた後で、インデックス・サーバ4030が活性化され、ログ済みのXQuery4020が、インデックス・テーブル4035を作るために解析される。このプロセスは、データベース4015内に格納されたデータのスキーマ・グラフ表現、および、ログ済みのXQuery4020によって表現されるその既存のビューを構築し、これらのビューに関連する様々な頻度テーブルを作り、データベース内の検索可能なキーワードを識別し、1つまたは複数のコンテキストツリーの決定およびコンテキストツリーそれぞれに対応するXQueryを構築し、最後に、これらのキーワードおよびXQueryのインデックス・テーブル4035への記録を、後で迅速に取り出せるように行うことを含む。

10

【0045】

インデックス・テーブル4035の作成プロセスが完了した後は、システム4000では、キーワード検索クライアント4025で呼び出されたキーワード検索操作を行う準備ができています。ユーザの入力した検索キーワード4026は、インデックス・サーバ4030へと送信され、ここで、そのルック・アップがインデックス・テーブル4035と突き合わせて行われ、1つまたは複数のXQuery4031が取り出され、検索キーワードにとっての関連性の高さに従って適切にランク付けされて、ユーザへと提示される。ユーザがXQuery4027をリストから選択すると、そのXQuery4027は、キーワード検索クライアント4025によってデータ・サーバ4010へと送信され、これは適切なデータ4011で応答する。

20

【0046】

1つまたは複数の階層的なデータソースを含むキーワード検索の方法2000を、図2のフローチャートで要約している。方法2000は、好ましくは、インデックス・サーバ4030のコンピュータ上で実行される。方法2000は、ステップ2005で始まり、ここでスキーマ・グラフ内のヒット・ノードが識別される。XMLの環境では、ヒット・ノードがスキーマ・グラフ内で生じる可能性のある仕方が2つある。すなわち、(i)その要素名が検索キーワードのうちの1つを含む可能性があり、または(ii)それが表現する1つまたは複数のXMLノードが検索キーワードのうちの1つを含む可能性がある。ステップ2005に続いて、ステップ2010では、スキーマ・グラフ内のヒット・ノードおよびコンテキスト・ノードで表現されるデータソース内のノードをそれぞれが含む、スキーマ・グラフ内のコンテキストツリーが識別される。最後にステップ2015で、識別されたコンテキストツリーがXQueryへと変換され、ランク付けされたリストとしてユーザに提示される。

30

【0047】

図2のステップ2010の表すコンテキストツリーを識別するための方法の説明を、ここで詳細に行う。まず、スキーマ・グラフ内に単一のヒット・ノードがある特別な場合に対する方法を提示し、その後、1つより多くのヒット・ノードを含む場合を扱えるより一般的な方法が続く。どちらの方法も、2フェーズで動作する。その第1は、スキーマ・グラフのヒット・ノードからのボトム・アップ走査(traversal)であって、その親および祖先のうちのどれがコンテキスト・ノードであるかを判定し、第2のフェーズは、これからトップ・ダウン的に進んで、その子孫のうちのどれが同様にコンテキスト・ノードであるかを判定する。次いで、ヒット・ノードの、コンテキスト・ノードと判定された最上位の祖先は、キーワード検索操作の結果としてユーザへと提示されるコンテキストツリーのルート・ノードに相当する。スキーマ・グラフ内のノードがコンテキスト・ノード

40

50

ドであるかどうかを判定する目的で、好ましくは、少なくとも生起頻度テーブルおよび共起頻度テーブルが維持される。前者では、スキーマ・グラフ内の各ノードがログ済みの問い合わせまたはデータ・ビューの中で生起する頻度が記録されるのに対し、後者では、スキーマ・グラフ内のノードの対が、同じログ済みの問い合わせまたはデータ・ビューの中で共起する頻度が記録される。スキーマ・グラフが、新しいノードを含む新しい問い合わせまたはデータ・ビューで更新されると、この新しいノードに相当する新しいエントリが生起頻度テーブルに追加され、それぞれには、そのノードが新しく、それまでに観測されたことがないことを示す頻度の初期値 1 が与えられる。同様に、2つの新しいノードまたは新しいノードと既存のノードを含む、新しい問い合わせからのノード対それぞれについては、新しいエントリが共起頻度テーブルに追加され、頻度の初期値 1 が与えられるのに対して、新しい問い合わせの中にない、新しいノードと既存のノードを含むノード対それぞれについては、新しいエントリは共起頻度テーブルに追加されるが、頻度の初期値 0 が与えられる。

10

【0048】

スキーマ・グラフが走査される際、ヒット・ノードの生起数のもとの、ノードごとの生起確率が計算される。こうした条件付き確率の値の計算または近似は、それまでにログの行われた問い合わせまたはデータ・ビューが原因で頻度テーブルに保存された値から行われ、あるノードがコンテキスト・ノードであるかどうかの判定に使用される。

【0049】

以下は、単一のヒット・ノードがある特別な場合に対する第 1 の方法の説明である。このヒット・ノードを X で表すことにする。

20

【0050】

第 1 のフェーズでは、スキーマ・グラフを通してのボトム・アップ走査はノード X を始めに行われる。X の祖先 Y_i のそれぞれを順番に考え、その X の生起のもとの生起確率を計算する。

【0051】

【数 1】

$$\Pr[Y_i | X] = \frac{\Pr[Y_i \wedge X]}{\Pr[X]}$$

30

$$\approx \frac{\text{freq}(Y_i, X)}{\text{freq}(X)}$$

式 1

【0052】

ここで、確率の値は、生起頻度 $\text{freq}(X)$ および共起頻度 $\text{freq}(Y_i, X)$ で近似されている。後者は、X と Y_i が共起する頻度を表し、ここで Y_i は X の祖先である。どちらも、前に述べた生起頻度テーブルおよび共起頻度テーブルから直接に得られる。X の祖先ノードに関して計算されたこれらの確率の値から、X のもとで、特定の祖先 Y_i がルート・ノードである確率を決定することが可能である。 Z_1, \dots, Z_n が Y_i の親ノードを表すとすると、

40

$$\Pr[Y_i \text{ root} | X] = \Pr[\neg Z_1 \dots \neg Z_n | Y_i, X] \quad \text{式 2}$$

2

となる。

【0053】

すなわち、X のもとで、 Y_i がルートである確率は、 Y_i が存在しその親のどれもが存在しない確率である。式 2 の右辺を展開すると、

$$\begin{aligned} \Pr[Y_i \text{ root} | X] &= \Pr[\neg Z_1 \dots \neg Z_n | Y_i, X] \Pr[Y_i | X] \\ &= (1 - \Pr[Z_1 \dots Z_n | Y_i, X]) \Pr \end{aligned}$$

50

[Y_i | X]
3 式

が得られる。

【 0 0 5 4 】

Z_1, \dots, Z_n は、 Y_i が与えられると互いに排他的であることから (Y_i は、実際のどのような階層的なデータ構造でも、多くとも1つの親をもつことができる)、

【 0 0 5 5 】

【 数 2 】

$$\begin{aligned} \Pr[Y_i \text{ root} | X] &= \left(1 - \sum_j \Pr[Z_j | Y_i \wedge X] \right) \Pr[Y_i | X] & 10 \\ &= \Pr[Y_i | X] - \sum_j \Pr[Z_j | Y_i \wedge X] \Pr[Y_i | X] \\ &= \Pr[Y_i | X] - \sum_j \Pr[Z_j \wedge Y_i | X] & \text{式 4} \end{aligned}$$

【 0 0 5 6 】

しかし、 Z_j と X の間には多くとも1つの有向パスしかない (スキーマ・グラフの1つの特徴) ことから、このパスは Y_i を含まなければならないことになり、これより

$$\Pr[Z_j | Y_i \wedge X] = \Pr[Z_j | X] \quad \text{式 5}$$

$$\Pr[Z_j | Y_i | X] = \Pr[Z_j | X] \quad \text{式 6}$$

【 0 0 5 7 】

【 数 3 】

$$\Leftrightarrow \Pr[Y_i \text{ root} | X] = \Pr[Y_i | X] - \sum_j \Pr[Z_j | X] \quad \text{式 7}$$

【 0 0 5 8 】

となる。

【 0 0 5 9 】

好ましい一実装形態では、いくつもの代替のコンテキストツリーがキーワード検索操作の結果としてユーザへと、関連づけられている確率 $\Pr[Y_i \text{ root} | X]$ が 0 より大きい X の祖先ノード Y_i ごとに1つ、返される。これらの代替のコンテキストツリーは、それぞれ、関連づけられている確率 $\Pr[Y_i \text{ root} | X]$ が割り当てられ、このスコアに従って、最高から最低へとソートされている。スコアがより高いコンテキストツリーは、スコアがより低いものよりもユーザにとって興味があると考えられる。一代替実装形態では、スコアが最高のコンテキストツリーだけが、キーワード検索操作の結果としてユーザへと返される。

【 0 0 6 0 】

コンテキストツリーのルート・ノードの役割のできる (すなわち、その $\Pr[Y_i \text{ root} | X] > 0$) 祖先ノード Y_i それぞれについて、第2のフェーズ、すなわち Y_i からのトップ・ダウン走査が、その子孫 (ヒット・ノード X を除く) のうちのどれがコンテキスト・ノードであるかを判定するために行われる。このフェーズ中に訪問した親ノード P_j それぞれについて、解析が、その子のうちのどれがコンテキスト・ノードであるかを判定するために行われる。コンテキスト・ノードであると判定された子ノードそれぞれについて、その子が順番にトップ・ダウン的に解析されて、そのうちのコンテキスト・ノードが識別される。

【 0 0 6 1 】

親ノード P_j の解析には、図 3 A および 3 B に示すように、2つの別々の想定状況がある。第1のものは、図 3 A に示す特別な場合であり、ここでは、 P_j はルート・ノード Y

40

30

50

i 3005からヒット・ノード X 3020へのパスに沿って存在する。これには、 $P_j = Y_i$ の場合が含まれるが、 $P_j = X$ の場合は除外される。この想定状況では、 P_j から X へのパスに沿って存在する、 P_j 3010の少なくとも1つの子ノード、すなわちこの場合では子ノード C_i 3015が、コンテキスト・ノードとして識別されなければならない。より一般的な、第2の想定状況では、図3に示すように残りの場合すべてを包括しているが、親ノード P_j 3030は、ルート・ノード Y_i からヒット・ノード X 3035への有向パスに沿っては存在せず、このため、 P_j のどれかの子ノードをコンテキスト・ノードとして識別することは義務的ではない。この第2の想定状況を扱うためのアルゴリズムを最初に提示することにする。

【0062】

10

所与のヒット・ノード X 、および、コンテキストツリーのルート・ノードの役割をする特定のノード Y_i に関して、親ノード P_j のある子ノード C_k をコンテキスト・ノードとして識別すべきかどうかを選ぶことは、一般に、 Y_i から X への有向パスに沿ったノードすべての存在のもとで、 C_k が生起する確率の関数

$$\Pr [C_k | X \dots Y_i \text{ root } \dots P_j]$$

である。

【0063】

この確率の評価または推定は、前にふれた生起および共起の頻度テーブルだけでは可能ではないため、何らかの形の単純化または近似が必要である。採用することが好ましいそのような単純化の1つは、上の確率の式の中の Y_i から C_k までのノード以外のノードすべての影響を無視することであり、次の式

20

$$\Pr [C_k | Y_i \text{ root } \dots P_j]$$

ここで

【0064】

【数4】

$$\Pr[C_k | Y_i, \text{root} \wedge \dots \wedge P_j] = \Pr[C_k | Y_i, \text{root} \wedge P_j]$$

$$= \frac{\Pr[C_k \wedge Y_i, \text{root} \wedge P_j]}{\Pr[Y_i, \text{root} \wedge P_j]}$$

30

$$= \frac{\Pr[C_k \wedge Y_i, \text{root}]}{\Pr[Y_i, \text{root} \wedge P_j]} \tag{式8}$$

【0065】

Z_l 、 $l = 1, \dots, p$ が Y_i の親ノードを表すとすると、式8の右辺は、

【0066】

【数5】

40

$$\frac{\Pr[C_k \wedge Y_i] - \sum_{l=1}^p \Pr[C_k \wedge Z_l]}{\Pr[P_j \wedge Y_i] - \sum_{l=1}^p \Pr[P_j \wedge Z_l]} \approx \frac{\text{freq}(Y_i, C_k) - \sum_{l=1}^p \text{freq}(Z_l, C_k)}{\text{freq}(Y_i, P_j) - \sum_{l=1}^p \text{freq}(Z_l, P_j)} \tag{式9}$$

【0067】

と展開することができる。

【0068】

しかし、上の式は、個々の子ノード C_k それぞれを孤立して扱っているのにすぎない。 C_k が(Y_i rootのもとで)互いに独立ない限りは、それらの同時確率を考える必

50

要がある。しかし、これには、ノードからなる大量の組合せの同時生起頻度を格納する頻度テーブルを維持することが必要であるが、そのうちの多くは、ほとんど観測されず、そのため、その同時確率の信頼できる推定をその同時生起頻度から行うことは可能ではないはずである。他方、 $(Y_i \text{ root})$ のもとで C_k の間の独立を仮定すると、子ノード C_k が、その個々の $(Y_i \text{ root})$ のもとでの生起確率が低い場合に、どれも選択されないなどの望ましくない結果をもたらされる。

【0069】

C_k の間の独立の過程の望ましくない影響を避ける一方で、同時に、大量の同時生起頻度の値を維持する必要性を避けるために、図5のフローチャートで示すヒューリスティックな方法5000を、子ノードをコンテキスト・ノードとして選択するのに使用することができる。方法5000は、好ましくは、サーバ4030上の方法2000内部のサブプログラムとして動作する。

【0070】

方法5000は、ステップ5005で始まり、ここで、 $Q_k = \text{Pr}[C_k | Y_i \text{ root} \dots P_j]$ で表される、ルート・ノード Y_i のもとの子ノード C_k それぞれの生起確率の計算が、式8および式9を用いて行われる。次のステップ5010で、確率 Q_k が子ノード C_k すべてにわたって総計されるが、その総和は T で表される。方法5000は継続され、ステップ5015で確率の値が最も大きいノード C_k がコンテキスト・ノードとして選択される。確率の値が同じで最大である子ノードが1つより多くある場合は、そのようなノードはすべてコンテキスト・ノードとして選択される。次いで、それまでにコンテキスト・ノードとして選択された子ノードすべての確率の総和 S の計算が、ステップ5020で行われる。次いで、実行は判断ステップ5025へと進み、このポイントで、子ノード C_k がすべてコンテキスト・ノードとして選択されている場合には、方法はステップ5040で終了する。しかし、まだコンテキスト・ノードとして選択されていない1つまたは複数の子ノード C_k がある場合は、方法5000はもう1つの判断ステップ5030へと続く。ステップ5030で、 $S > T/2$ かどうかを確かめるために検査が行われ、そうである場合、方法はここでもステップ5040で終了する。 $S < T/2$ である場合、実行はステップ5035へと進む。ここでは、現在まだコンテキスト・ノードとして選択されていない子ノード C_k のリストが検査されて、そのうちで確率が最大のものが識別される。それがコンテキスト・ノードとして選択され、方法5000は、ステップ5020へと戻って、上で論じたような処理をさらに行う。

【0071】

方法5000には、いくつかの望ましい特徴がある。すなわち、

- ・ログ済みの問い合わせまたはデータ・ビューに共通部分があって頻度が十分に高い(すなわち、比較的多数の子ノードが共通にある)場合、方法5000はその共通部分をコンテキスト・ノードとして返す傾向がある。これからは、十分に大きな共通部分は(ヒット・ノードに関する)十分なコンテキスト情報を担う傾向があるため、容認できる結果をもたらされそうである。

【0072】

- ・ログ済みの問い合わせまたはデータ・ビューに、共通に比較的少数の子ノードしかなかった場合、結果となるコンテキスト・ノードの組は、その共通部分だけでなくさらに別のノードも含む傾向がある。発明者の行った実験では、結果となるコンテキストの子ノードの組は、頻度が最も高いログ済みの問い合わせまたはデータ・ビューのそれを反映する傾向がある。このことは、共通部分だけでは十分なコンテキスト情報を含むことはありそうにないであろうから、重要である。

【0073】

- ・確率の値が等しく最大である子ノードをまとめて含めることが原因で、方法5000は、ノードのコンテキスト・ノードとしての識別を、少なくではなく多く行う方へのバイアスをもつ。ログ済みの問い合わせの中にある子ノードの組が互いに排他的であり等しい頻度で生じる場合には、この方法では、子ノードすべてがコンテキスト・ノードとして識

10

20

30

40

50

別される。

【0074】

方法5000では、親ノードがコンテキスト・ノードと同じであるときには、その子のうちの1つまたは複数も、必ずコンテキスト・ノードとして識別される。このことは、親ノードが、まったく子がない状態で生じる（すなわち、葉ノードとして生じる）ログ済みの問い合わせまたはデータ・ビューが数多くある場合には、望ましくない。直観的には、これが十分しばしば生じる場合には、その親ノードだけを、まったく子がないコンテキスト・ノードとして識別して、頻繁に観察される振る舞いを反映させるべきである。

【0075】

この問題点を改善するために、好ましい一実装形態では、インデックス・サーバ4030の生成および保存する追加の葉共起頻度テーブルを利用する。このテーブルでは、ノード P_j が、葉ノードとして過去のログ済みの問い合わせおよびデータ・ビュー内でその祖先 Y_i と共起している頻度を、スキーマ・グラフ内に子のないノード P_j を除き、そのようなノード P_j と Y_i のあらゆる可能な対に対して保存する。次いで、この新しい頻度テーブルを使用して、 P_j およびあるルート・ノード Y_i のもとで、ノード P_j が葉ノードとして生起する確率

【0076】

【数6】

$$\Pr[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] = \frac{\Pr[P_j \text{ leaf} \wedge Y_i \text{ root}]}{\Pr[Y_i \text{ root} \wedge P_j]}$$

10

20

$$\approx \frac{\text{freq}(Y_i, P_j \text{ leaf}) - \sum_{l=1}^p \text{freq}(Z_l, P_j \text{ leaf})}{\text{freq}(Y_i, P_j) - \sum_{l=1}^p \text{freq}(Z_l, P_j)}$$

式10

【0077】

が推定される。ここで、 Z_l 、 $l = 1, \dots, p$ は、前に定義したように Y_i の親ノードを表し、 $\text{freq}(Y_i, P_j \text{ leaf})$ および $\text{freq}(Z_l, P_j \text{ leaf})$ は、新しい葉共起頻度テーブルから得られる共起頻度の値である。

30

【0078】

確率 $\Pr[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j]$ は、好ましくは、図5に与える方法5000に先立つ追加の判断ステップで決定され、 P_j のどの子ノードがコンテキスト・ノードであるかが識別される。 $\Pr[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j]$ が0.5より小さい場合は、 P_j の子ノードはコンテキスト・ノードとしては選択されず、そうでない場合は、方法5000が行われて、どの子ノードがコンテキスト・ノードであるかが識別される。

【0079】

代替実施形態も可能であり、そこでは、フローチャートを図6に与えている代替の方法6000を使用して、1組の子ノード C_k 、 $k = 1, \dots, m$ のうちでコンテキスト・ノードが選択される。方法6000は、これもインデックス・サーバ4030で行われるが、ステップ6001で始まり、そこでは、概念上、架空の子ノード C_0 が作成され、実際の子ノード C_1, \dots, C_m からなるリストへと追加され、確率値 $\Pr[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j]$ が式10を用いて割り当てられる。次のステップ6005で、実際の子ノード C_k には、その通常確率値 $Q_k = \Pr[C_k | Y_i \text{ root} \wedge P_j]$ が式8および式9を用いて割り当てられる。次いで、方法6000は継続され、ステップ6006で、方法5000をステップ5010で呼び出して（ステップ5005を飛ばして）子ノード C_0, \dots, C_m のうちで1組のコンテキスト・ノードを選択する。方法5000を抜けると、方法6000は再開され、判断ステップ6010で検査が行わ

40

50

れて、架空の子ノード C_0 がコンテキスト・ノードとして選択されたかどうか判定される。そうである場合、実行は継続され、ステップ 6020 で、 C_0 がコンテキスト・ノードとして除外される。方法 6000 は、その後、ステップ 6015 で終了する。検査が 6010 で失敗した場合、方法 6000 は、直接に終了ステップ 6015 へと進む。

【0080】

P_j の子ノードのうちのどれもコンテキスト・ノードではないという可能性を組み込むための代替の方法 6000 の背後にある考え方は、第 1 の方法と本質的に同じである。すなわち、 $\Pr[P_j \text{ leaf} | Y_i \text{ root } P_j]$ が十分に大きいときである。しかし、コンテキスト・ノードの最終的なセットに対する $\Pr[P_j \text{ leaf} | Y_i \text{ root } P_j]$ の影響は、この代替のアプローチでは、よりゆるやかであり、一般に、第 1 のアプローチでの急激なオン・オフの切り替えよりも好適である。

10

【0081】

親ノード P_j がルート・ノード Y_i からヒット・ノード X への有向パスに沿って存在する特別な想定状況では、 P_j から X へのパスに沿って存在する P_j の子ノードがコンテキスト・ノードとして識別されるように、特別な考慮がなされなければならない。一般性を失うことなく、この子ノードを、図 3 に項目 3015 として示す C_1 とする。前に提示した一般的な想定状況用の方法 5000 は、(たとえば、 C_1 の生起確率を他の子ノードすべてのそれを上回るように、ステップ 5015 の前に、吊り上げておくことによって) 変更が可能であるが、そのようなアプローチは正しい結果を生じない可能性がある。その理由は、説明した方法 5000 は、ルート Y_i および親 P_j のもとでコンテキスト・ノードとして最も頻繁に生起する子ノードの組を選択するように工夫されているためである。この組が自然に C_1 を含まない場合、それは基本的に、 C_1 がその組の中のノードと関係がないことを意味している。 C_1 を無理に含めても、ただ、共通するところが少ししかない子ノードの組ができてしまい、 C_1 に対する(またその後 X に対する)コンテキストは少ししか得られないはずである。

20

【0082】

方法 5000 を変更する代わりに、好ましくは、異なるが手続き的にはいくぶん似ている、図 7 に示す方法 7000 が、別の実装形態では採用される。新しい 7000 と以前の 5000 の間の違いは、使用される確率の独立性の仮定にある。 P_j が Y_i から X への有向パスに沿って存在しない、一般的な場合に行った第 1 の単純化が、

30

$$\Pr[C_k | X \dots Y_i \text{ root } \dots P_j]$$

は、 Y_i から P_j までのノード以外のノードとは独立であるという仮定であったことを想起されたい。 P_j の 1 つの子ノード C_1 が X から P_j へのパスに沿って存在するという現在の想定状況では、 C_k が (C_1 を含めて) P_j から X までのノードとは独立であると仮定することは、これらがヒット・キーワード X を C_k にリンクする X の必要な祖先であることから、賢明ではあるまい。問題を扱いやすくしておくには単純化がいくつか必要であるため、 C_k とその P_j より上位のルート・ノード Y_i に向かう祖先の間の確率の独立性を仮定することはよい選択肢ということになる。この仮定のもとで、興味の対象となる確率は、

$$\Pr[C_k | X \dots C_1 P_j] \quad k \geq 1$$

40

である。

【0083】

再び、 X と P_j をリンクする有向パスが多くとも 1 つあるため、上の式は

【0084】

【数 7】

$$\Pr[C_k | X \wedge P_j] = \frac{\Pr[C_k \wedge X \wedge P_j]}{\Pr[X \wedge P_j]}$$

式 11

【0085】

に等しい。

50

【0086】

式11の右辺の分子は、2つではなく3つのノードを含んでいる以上、これまでにふれた生起および共起の頻度テーブルからは得ることができない。したがって、ノードの3つ組(3-tuple)の間の、追加の同時生起頻度テーブルが必要である。さいわい、こうした3つ組は、(ノードの任意のどの対でもではなく)親子ノードC_kとP_jからなる対を含んでおり、各ノードC_kには、実際には、少数の親しかいないため、新しい同時生起頻度テーブルは、ノードの対を含む共起頻度テーブルよりも、わずかに大きいのにすぎないはずである。

【0087】

新しい同時生起頻度テーブルにより、Pr[C_k | X P_j]は

10

【0088】

【数8】

$$Pr[C_k | X \wedge P_j] \approx \frac{freq(C_k, P_j, X)}{freq(P_j, X)} \tag{式12}$$

【0089】

と推定することができ、ここで、freq(C_k, P_j, X)は、ノードC_k、P_j、およびXの間の同時生起頻度を表し、P_jはC_kの親でXの祖先であり、C_kはXでもXの祖先でもない。

【0090】

20

C₁とともにコンテキスト・ノードとして含めるべきC₁の兄弟(siblings)からなる組を決定するための方法7000は、すでに説明した方法5000に非常に似ている。方法7000は、ステップ7001で始まり、ここでは、Q_k = Pr[C_k | X P_j]で表される、親ノードP_jおよびヒット・ノードXのもとでの各子ノードC_k C₁の生起確率の計算が式12を用いて行われる。次のステップ7005で、確率Q_kが子ノードC_k C₁すべてにわたって総計されるが、その総和はTで表される。方法7000は継続され、ステップ7010で、ノードC₁がコンテキスト・ノードとして選択され、次いでその後ステップ7015で、確率の値が最大のC_k C₁もコンテキスト・ノードとして選択される。確率の値が最大で同じ子ノードが1つより多く存在する場合は、そのようなノードはすべてコンテキスト・ノードとして選択される。次いで、C₁を除きそれまでにコンテキスト・ノードとして選択した子ノードすべての確率の総和の計算を、ステップ7020で行い、その総和はSで表される。次いで、実行は、判断ステップ7025へと進み、このポイントで、子ノードC_kがすべてコンテキスト・ノードとして選択されている場合、方法7000は、ステップ7040で終了する。しかし、1つまたは複数の子ノードC_kがまだコンテキスト・ノードとして選択されていない場合は、この方法は、別の判断ステップ7030へと続く。ステップ7030で、検査を行ってS > T/2であるかどうかを確かめ、そうである場合は、方法7000は、ここでも、ステップ7040で終了する。S < T/2である場合は、実行はステップ7035へと進む。ここで、現在まだコンテキスト・ノードとして選択されていない子ノードC_k C₁からなるリストが検査されて、それらのうちで確率の値が最大のものが識別される。これらは、コンテキスト・ノードとして選択され、方法は、処理をさらに行うためにステップ7020へと戻る。

30

40

【0091】

C₁の兄弟が解に含まれていない場合を見越して、方法7000にいくつかの変更が必要である。これは、ノードP_jが、過去のログ済みの問い合わせおよびデータ・ビューにP_jの子ノードが1つしかない(P_jからXへのパスに沿うC₁)のようなその子孫Xのうちの1つと共起する頻度の格納されている単独子共起頻度テーブルを導入することによって達成される。次いで、この頻度テーブルを使用して、親P_jおよびヒット・ノードXのもとでC₁に兄弟がない確率

【0092】

50

【数9】

$$\Pr[C_1, \text{no sibling} | P_j \wedge X] = \Pr[C_1 \wedge \neg C_k \forall k \neq 1 | P_j \wedge X]$$

$$= \Pr[P_j \text{ has 1 child} | P_j \wedge X]$$

$$= \frac{\Pr[P_j \text{ has 1 child} \wedge P_j \wedge X]}{\Pr[P_j \wedge X]}$$

$$= \frac{\Pr[P_j \text{ has 1 child} \wedge X]}{\Pr[P_j \wedge X]}$$

$$\approx \frac{\text{freq}(P_j \text{ has 1 child}, X)}{\text{freq}(P_j, X)}$$

式13

10

【0093】

が推定され、ここで、 $\text{freq}(P_j \text{ has 1 child}, X)$ は、 P_j がその子孫 X と共起し、 P_j に単一の子ノード(C_1)がある頻度を表しており、新しい頻度テーブルから得られる。

20

【0094】

ある実装形態では、確率 $\Pr[C_1 \text{ no sibling} | P_j \wedge X]$ を、図7に与えている方法7000に先立つ追加の判断ステップで使用して、 P_j のどの子ノードがコンテキスト・ノードであるかが識別される。 $\Pr[C_1 \text{ no sibling} | P_j \wedge X]$ が0.5より小さい場合は、 C_1 以外の P_j の子ノードはコンテキスト・ノードとして選択されず、そうでない場合は、方法7000が行われてどの子ノードがコンテキスト・ノードであるかが識別される。

【0095】

代替実施形態も可能である。ここでは、フローチャートを図8に与えている代替の方法8000を使用して、1組の子ノード C_k 、 $k=1, \dots, m$ のうちでコンテキスト・ノードが選択される。方法8000は、ステップ8001で始まり、そこでは、概念上、架空の子ノード C_0 が作成され、実際の子ノード C_1, \dots, C_m からなるリストへと追加され、確率値 $Q_0 = \Pr[C_1 \text{ no sibling} | P_j \wedge X]$ が式13を用いて割り当てられる。ステップ8005で、 C_1 を除く実際の子ノード C_k には、その通常確率値 $Q_k = \Pr[C_k | X \wedge P_j]$ が式11を用いて割り当てられる。次いで、方法8000は継続され、ステップ8006で、方法7000をステップ7005で呼び出して(ステップ7001を飛ばして)子ノード C_0, \dots, C_m のうちで1組のコンテキスト・ノードを選択する。方法7000を抜けると、方法8000は再開され、判断ステップ8010で検査が行われて、架空の子ノード C_0 がコンテキスト・ノードとして選択されたかどうか判定される。そうである場合、実行は継続され、ステップ8020で、 C_0 がコンテキスト・ノードとして除外される。方法8000は、その後、ステップ8015で終了する。検査が8010で失敗した場合、方法8000は、直接に終了ステップ8015へと進む。

30

40

【0096】

これまでの議論で、コンテキスト・ノードである1組の子ノードから決定するための2つの異なる方法6000および8000を説明している。好ましくは、後者は、親ノード P_j がルート・ノード Y_i からヒット要素 X への有向パスに沿って存在する想定状況で適用されるのに対し、前者は他の親ノードすべてに対して使用される。一代替実装形態では、第1の方法6000は、 P_j が Y_i から X へのパスに沿って存在する場合でも使用される。これから C_1 を含む1組のコンテキストの子ノードができる場合は、その組が採用さ

50

れるが、そうでない場合は、その組は捨てられ、第2の方法8000が、コンテキストの子ノードからなる新しい組を決定するために適用される。この第1の方法の有利さの背後にある理論的根拠は、そこで計算される確率の値が条件としているのが、ルート要素 Y_i であり、ヒット・ノード X ではないことである。本発明者の行った試験は、データ・ビューのルート要素が、どのようなノードがそのビュー内にあるかのよりよい指標であることを示唆するように思われる。

【0097】

本明細書で開示するキーワード検索システム4000は、学習システムの一形式である。ログ済みの問い合わせおよび既存のデータ・ビューの組、これはトレーニングの例に似ているが、これからシステムは、データの新しいビューを合成することができる。ログ済みの問い合わせまたはデータ・ビューの中にパターンが存在する場合、そのパターンは頻度テーブルの中に反映され、これが今度はシステム4000の振る舞いに影響することになる。どのような学習システムにとっても望ましい特徴は、その学習システムが、関係があるがまだ見たことのない問題を扱うときに、1組の問題から学習したパターンを使用してそのパフォーマンスを改善することができる何らかの形の一般化を行う能力である。階層的な環境で重要な一般化の1つの側面は、データの特定の下位構造の生起パターンを観察し、それを他の類似または同一の下位構造へと一般化する能力である。

【0098】

図9に示すデータ構造9000を考えるが、この中には、2つの同じ従業員(Employee)下位構造9010および9030が(点線の曲線で囲まれている)、1つは管理職(Manager)9005の下に、もう1つはプロジェクト・メンバー(Project Member)9025の下にある。ログ済みの問い合わせまたはデータ・ビューすべての中で、第1のEmployeeサブツリー内の下位要素の名(FirstName)9015および姓(LastName)9020が、ともに現れることがすでに観察されているのに対して、第2の従業員(Employee)サブツリー9030を含む問い合わせまたはデータ・ビューはまだ何も観察されていないとする。さらに、従業員の名前に対するキーワード検索操作が呼び出されて、「ヒット」が第2の従業員(Employee)サブツリー9030の名(FirstName)下位要素9035の中で見付かり、9035がヒット・ノードになっているとする。たとえ、この下位要素が存在する問い合わせまたはビューの例に出遭っていないとしても、第1の従業員(Employee)サブツリー9010に関して観察された生起パターンから、第2の従業員(Employee)サブツリー9030内の下位要素の姓(LastName)9040をコンテキスト・ノードとして識別すべきであることは、直観的に明らかである。

【0099】

そのような一般化の能力は、XMLデータの作業を行うときには、同一のデータ下位構造が、しばしば、データ階層内のいくつかの場所に存在するため(たとえば、参照されるスキーマ要素を使用する結果として)、特に重要である。そのようなものは、確率平均化(probability averaging)を通して実現することができる。確率平均化の動作は、スキーマ・グラフ内で名前またはIDまたはラベルが同一のノードの生起確率を適切に平均することによって行われる。ここで、確率平均化の適用の説明を、まずコンテキストツリーの構築の第1のトップ・ダウンのフェーズについて、そうしてその後で第2のボトム・アップのフェーズについて行う。

【0100】

第1のフェーズの動作は、確率の値 $Pr[Y_i | X]$ を利用していることを想起されたい。ここで、 Y_i はヒット・ノード X の祖先である。確率平均化を容易にするために、 $Pr[Y_i | X]$ を、好ましくは、次のように漸化式(incremental form)へと定式化し直す。すなわち、 W を、 Y_i から X への唯一の有向パスに沿って存在する、 Y_i の子とする。すると、 $Pr[Y_i | X]$ を

【0101】

10

20

30

40

【数10】

$$\begin{aligned}
 \Pr[Y_i | X] &= \frac{\Pr[Y_i \wedge X]}{\Pr[X]} \\
 &= \frac{\Pr[Y_i \wedge W \wedge X]}{\Pr[X]} && (Y_i \text{ から } X \text{ へのパスは } W \text{ を含む}) \\
 &= \frac{\Pr[Y_i | W \wedge X] \Pr[W \wedge X]}{\Pr[X]} \\
 &= \Pr[Y_i | W \wedge X] \Pr[W | X] && \text{式14}
 \end{aligned}$$

10

【0102】

のように書き直すことができる。

【0103】

すなわち、 $\Pr[Y_i | X]$ は、その子ノード W の確率の値、すなわち $\Pr[W | X]$ から漸進的に得ることができる。考え方は、手続きをヒット・ノード X で始め、上の式を利用して、次々に上位の祖先ノードに対する確率の値を得ることである。各ステップで、そのときに確率平均化の方法が、式14の右辺の第1項に適用される。たとえば、 $\Pr'[B | X]$ が、あるノード B の確率平均化の結果としての変更後の確率の値を表すとすると、 $\Pr'[Y_i | X]$ を、次の再帰式 (recursive formulae) で定義することができる。すなわち

20

$$\Pr'[X | X] = 1 \quad \text{式15}$$

【0104】

【数11】

$$\Pr'[Y_i | X] = \begin{cases} 0 & \text{if } \Pr'[W | X] = 0 \\ \Pr_{mean}[Y_i | W \wedge X] \Pr'[W | X] & \text{otherwise} \end{cases} \quad \text{式16}$$

【0105】

ここで

30

【0106】

【数12】

$$\begin{aligned}
 \Pr_{mean}[Y_i | W \wedge X] &= \frac{\sum_k \Pr[Y_{ik} | W_k \wedge X_k] \Pr[W_k \wedge X_k]}{\sum_k \Pr[W_k \wedge X_k]} \\
 &= \frac{\sum_k \Pr[Y_{ik} \wedge W_k \wedge X_k]}{\sum_k \Pr[W_k \wedge X_k]} && \text{40} \\
 &= \frac{\sum_k \Pr[Y_{ik} \wedge X_k]}{\sum_k \Pr[W_k \wedge X_k]} && (Y_{ik} \text{ から } X_k \text{ までのパスは } W_{ik} \text{ を含む}) \\
 &\approx \frac{\sum_k \text{freq}(Y_{ik}, X_k)}{\sum_k \text{freq}(W_k, X_k)} && \text{式17}
 \end{aligned}$$

50

【 0 1 0 7 】

であり、これが表すのは、 (Y_i, X) に等しい (k のある値に対する) ノードの対 $(Y_{i k}, X_k)$ すべてにわたって、 X_0 および $Y_{i 0}$ (すなわち、 $k = 0$) をそれぞれ X および Y_i の別名 (aliases) として計算した、 W および X のもとでの Y_i の重み付け平均 (average or mean) 確率である。これらの等しい対 $(Y_{i k}, X_k)$ のそれぞれに対して、総和の中の項 W_k は、 $Y_{i k}$ から X_k への有向パスに沿って存在する、 $Y_{i k}$ のすぐ下の (immediate) 子を表す。

【 0 1 0 8 】

ノード対 $(Y_{i k}, X_k)$ がノード対 (Y_i, X) に等しいというのは、

(i) $Y_{i k}$ の名前またはラベルまたは ID が Y_i と同じであり、 X_k の名前、ラベル、または ID が X と同じである、

(ii) $Y_{i k}$ と X_k の間に直接の祖先 - 子孫関係があり、 Y_i と X の間でも同様である、

(iii) $Y_{i k}$ から X_k への有向パスに沿うノード W_k それぞれに対して、 (W_k, X) が (W, X) に等しく、 (Y_i, W_k) が (Y_i, W) に等しいような Y_i から X への有向パスに沿う対応するノード W がなければならない、

(iv) Y_i と $Y_{i k}$ にちょうど同じ数の親があり、 Y_i の親 Z_j それぞれに対して、 $(Z_{k j}, Y_{i k})$ および (Z_j, Y_i) が上の条件 (i) から (iii) を満たすような $Y_{i k}$ の親 $Z_{k j}$ がある

場合である。

【 0 1 0 9 】

すると、確率平均化が原因で変更された、 X のもとで Y_i がルートである確率は、

【 0 1 1 0 】

【 数 1 3 】

$$\Pr[Y_i \text{ root} | X] = \Pr[Y_i | X] - \sum_j \Pr[Z_j | X] \quad \text{式18}$$

【 0 1 1 1 】

ここで、式 16 から、 Y_i を Z_j で、 W を Y_i で置き換えることによって得られるように、

$$\Pr[Z_j | X] = \Pr_{\text{mean}}[Z_j | Y_i, X] \Pr[Y_i | X] \quad \text{式19}$$

9

である、

式 17 の右辺の分母が 0 である場合、これはログ済みの問い合わせおよびデータ・ビューでノード対 (W_k, X_k) のうちのどれも観察されていないことを示すが、式 17、したがって式 19 および式 18 は不定 (undefined) となり、したがって、コンテキスト・ノードを識別するための何らかの代替方法が必要になる。好ましいアプローチは、 $\Pr_{\text{mean}}[Z_j | Y_i, X]$ の代替の定義を、次のように、 Z_j のヒット・ノード X からの距離によって行うことである。すなわち

【 0 1 1 2 】

10

20

30

40

【数14】

$Pr_{mean}[Z_j | Y_i \wedge X] =$

$$\begin{cases}
 \frac{\sum_k freq(Z_{jk}, X_k)}{\sum_k freq(Y_{ik}, X_k)} & \text{if } \sum_k freq(Y_{ik}, X_k) \neq 0 \\
 1 & \text{if } \sum_k freq(Y_{ik}, X_k) = 0, \text{ dist}(Z_j, X) \leq d_{max} \\
 0 & \text{if } \sum_k freq(Y_{ik}, X_k) = 0, \text{ dist}(Z_j, X) > d_{max}
 \end{cases}$$

式 20

10

【0113】

ここで、 d_{max} はある閾値の定数であり、 $dist(A, B)$ は、2つのノードAとBの間のスキーマ・グラフ内での距離であり、AとBの間のパスに沿うリンクの数として定義される。関連性の高い過去のログ済みの問い合わせおよびデータ・ビューがない場合は、2つのノードの間の距離によって、それらが互いにどう関係しているかのよい表示が与えられるが、これは、実際には、関係のあるデータは、普通、互いに近くに保存されているためである。

20

【0114】

ヒット・ノードXの祖先ノード Y_i がコンテキストツリーのルート・ノードである確率の計算を、確率平均化によって、祖先ノード Y_i すべてについて行うための方法10000のフローチャートを、図10に示している。方法10000は、ステップ10001で、 $Y_i = X$ 、したがって $Pr'[Y_i | X] = 1$ として始まる。次のステップ10005で、式19および式20を使用して、 $Pr'[Z_j | X]$ の計算が Y_i の親ノード Z_j それぞれに対して行われる。ステップ10005の後、ステップ10010では、 $Pr'[Y_i \text{ root} | X]$ の計算が式18に従って行われる。次いで、ステップ10025では、検査が、 Y_i の親ノードすべてが処理されているかどうかを判定するために行われる。そうでない場合、方法10000は、ステップ10015へと進み、ここで、 Y_i の親ノード Z_j が選択される。ステップ10020に到達すると、方法10000は、ステップ10005で(ステップ10001を飛ばして)再帰的に、ただし選択した親ノード Z_j が Y_i の役割を果たす状態で呼び出される。この呼び出しが戻ると、方法10000の現在の実行は再開され、ステップ10025へと戻って、親ノードがまだあるかの検査が行われる。親ノードがすべて処理されてしまうと、方法10000は、ステップ10030で終わる。

30

【0115】

確率平均化は、第2のトップ・ダウン走査のフェーズにも適用される。このフェーズで、親ノード P_j がルート・ノード Y_i からヒット・ノードXへの有向パスに沿って存在しない一般の場合には、確率平均化を第1のフェーズで使用したのと同じ方法で適用することができる。キーワード検索の結果にコンテキスト・ノードとして含めるための Y_i の子ノード C_k の選択は、確率

40

$$Pr[C_k | Y_i \text{ root } P_j]$$

に基づく。確率平均化によって、上の式は、平均確率

【0116】

【数 15】

$$\Pr_{mean}[C_k | Y_i \text{ root} \wedge P_j] = \frac{\sum_h \Pr[C_{kh} \wedge Y_{ih} \text{ root}]}{\sum_h \Pr[Y_{ih} \text{ root} \wedge P_{jh}]} \quad \text{式 21}$$

【0117】

によって置き換えられるが、ここで、 (Y_{ih}, C_{kh}) は (Y_i, C_k) に等しく、 (P_{jh}, C_{kh}) は (P_j, C_k) に等しく、 Y_{i0}, C_{k0} 、および P_{j0} (すなわち、 $h=0$) はそれぞれ、 Y_i, C_k 、および P_j の別名とする。 Z_j を Y_i の親、同様に Z_{jh} を Y_{ih} の対応する親とする。上の式を展開して

【0118】

【数 16】

$$\Pr_{mean}[C_k | Y_i \text{ root} \wedge P_j] = \frac{\sum_h \left\{ \Pr[C_{kh} \wedge Y_{ih}] - \sum_j \Pr[C_{kh} \wedge Z_{jh}] \right\}}{\sum_h \left\{ \Pr[P_{jh} \wedge Y_{ih}] - \sum_j \Pr[P_{jh} \wedge Z_{jh}] \right\}}$$

$$\approx \frac{\sum_h \left\{ \text{freq}(Y_{ih}, C_{kh}) - \sum_j \text{freq}(Z_{jh}, C_{kh}) \right\}}{\sum_h \left\{ \text{freq}(Y_{ih}, P_{jh}) - \sum_j \text{freq}(Z_{jh}, P_{jh}) \right\}} \quad \text{式 22}$$

【0119】

にすることができる。

【0120】

上の式が平均確率 $\Pr_{mean}[C_k | Y_i \text{ root} \wedge P_j]$ の正確な近似であるためには、右辺の分母が十分に大きくなければならない(たとえば、 $>$ ある正定数 f_{min})。そうでないときには、好ましい一実装形態において採用された好ましい改善方法を使用して、まず $\Pr_{mean}[C_k | Y_i \text{ root} \wedge P_j]$ が、 $\Pr_{mean}[C_k | Y_i \wedge P_j]$ で近似され、ここで、確率は、 $Y_i \text{ root} \wedge P_j$ でなく、 $Y_i \wedge P_j$ を条件としている。こうして

【0121】

【数 17】

$$\Pr_{mean}[C_k | Y_i \text{ root} \wedge P_j] \approx \Pr_{mean}[C_k | Y_i \wedge P_j]$$

$$= \frac{\sum_h \Pr[C_{kh} \wedge Y_{ih}]}{\sum_h \Pr[Y_{ih} \wedge P_{jh}]}$$

$$\approx \frac{\sum_h \text{freq}(Y_{ih}, C_{kh})}{\sum_h \text{freq}(Y_{ih}, P_{jh})} \quad \text{式 23}$$

【0122】

50

式 23 の右辺の分母がやはり十分に大きくない場合は、 $\Pr_{mean}[C_k | Y_i \text{ root } P_j]$ はさらに、 Y_i ではなく W を条件とする確率で近似され、ここで、 W は Y_i のすぐ下の子であり、 C_k の祖先である。すなわち

【 0 1 2 3 】

【 数 1 8 】

$$\Pr_{mean}[C_k | Y_i \text{ root } \wedge P_j] \approx \Pr_{mean}[C_k | W \wedge P_j]$$

$$= \frac{\sum_h \Pr[C_{kh} \wedge W_h]}{\sum_h \Pr[W_h \wedge P_{jh}]}$$

10

$$\approx \frac{\sum_h \text{freq}(W_h, C_{kh})}{\sum_h \text{freq}(W_h, P_{jh})}$$

式 24

【 0 1 2 4 】

この方法は繰り返して、さらに、右辺の分母に十分に大きな値が得られるまで、そうでない場合は、 W が C_k の親を表すまで行われる。後者の場合、 $\Pr_{mean}[C_k | Y_i \text{ root } P_j]$ への値の割り当ては、 C_k と Y_i の間の距離に基づいて

20

【 0 1 2 5 】

【 数 1 9 】

$$\Pr_{mean}[C_k | Y_i \text{ root } \wedge P_j] \approx \begin{cases} 1 & \text{if } \text{dist}(C_k, Y_i) \leq d_{\max} \\ 0 & \text{otherwise} \end{cases}$$

式 25

【 0 1 2 6 】

または、 C_k と ヒット・ノード X の間の距離に基づいて

【 0 1 2 7 】

【 数 2 0 】

$$\Pr_{mean}[C_k | Y_i \text{ root}] \approx \begin{cases} 1 & \text{if } \text{dist}(C_k, X) \leq d_{\max} \\ 0 & \text{otherwise} \end{cases}$$

式 26

【 0 1 2 8 】

のように行われる。

【 0 1 2 9 】

$\Pr_{mean}[C_k | Y_i \text{ root } P_j]$ の近似が、最終的に式 22、式 23、式 24、式 25、または式 26 で行われたのに応じて、ルート・ノード Y_i のもとで、親ノード P_j にコンテキストの子ノードがない平均確率の計算は、それぞれ、式 27、式 28、式 29、式 30、または式 31 を用いて行われる。

40

【 0 1 3 0 】

【数 2 1】

$$\Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] = \frac{\sum_h \Pr[P_{jh} \text{ leaf} \wedge Y_{ih} \text{ root}]}{\sum_h \Pr[Y_{ih} \text{ root} \wedge P_{jh}]}$$

$$\approx \frac{\sum_h \left\{ \text{freq}(Y_{ih}, P_{jh} \text{ leaf}) - \sum_j \text{freq}(Z_{jh}, P_{jh} \text{ leaf}) \right\}}{\sum_h \left\{ \text{freq}(Y_{ih}, P_{jh}) - \sum_j \text{freq}(Z_{jh}, P_{jh}) \right\}} \quad \text{式 27}$$

10

【0 1 3 1】

【数 2 2】

$$\Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] \approx \Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \wedge P_j]$$

$$= \frac{\sum_h \Pr[P_{jh} \text{ leaf} \wedge Y_{ih}]}{\sum_h \Pr[Y_{ih} \wedge P_{jh}]}$$

20

$$\approx \frac{\sum_h \text{freq}(Y_{ih}, P_{jh} \text{ leaf})}{\sum_h \text{freq}(Y_{ih}, P_{jh})} \quad \text{式 28}$$

【0 1 3 2】

【数 2 3】

$$\Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] \approx \Pr_{\text{mean}}[P_j \text{ leaf} | W \wedge P_j]$$

$$= \frac{\sum_h \Pr[P_{jh} \text{ leaf} \wedge W_h]}{\sum_h \Pr[W_h \wedge P_{jh}]}$$

30

$$\approx \frac{\sum_h \text{freq}(W_h, P_{jh} \text{ leaf})}{\sum_h \text{freq}(W_h, P_{jh})} \quad \text{式 29}$$

【0 1 3 3】

【数 2 4】

$$\Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] \approx \begin{cases} 0 & \text{if } \text{dist}(P, Y_i) + 1 \leq d_{\text{max}} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 30}$$

【0 1 3 4】

【数 2 5】

$$\Pr_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root} \wedge P_j] \approx \begin{cases} 0 & \text{if } \text{dist}(P, X) + 1 \leq d_{\text{max}} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 31}$$

【0 1 3 5】

50

ルート要素 Y_i のもとでの親ノード P_j に対するコンテキストの子ノードの決定を、 P_j が Y_i からヒット・ノード X への有向パスに沿って存在しない一般の場合について、確率平均化によって行うための好ましい手続きは、図 6 に示すものに類似しており、これを図 13 に示している。方法 13000 は、ステップ 13001 で始まり、そこでは、概念上、架空の子ノード C_0 が作成され、実際の子ノード C_1, \dots, C_m からなるリストへと追加され、確率値 $Q_0 = \text{Pr}_{\text{mean}}[P_j \text{ leaf} | Y_i \text{ root } P_j]$ が、式 27、式 28、式 29、式 30、または式 31 を用いて計算されて割り当てられ、次のステップ 13005 で、実際の子ノード C_k には、それに対応して、確率値 $Q_k = \text{Pr}_{\text{mean}}[C_k | Y_i \text{ root } P_j]$ が、それぞれ、式 22、式 23、式 24、式 25、または式 26 を用いて計算されて割り当てられる。どの場合でも、ステップ 13006 がステップ 13005 に続き、方法 5000 をステップ 5010 で呼び出して（ステップ 5005 を飛ばして）、子ノード C_0, \dots, C_m のうちで 1 組のコンテキスト・ノードを選択する。方法 5000 を抜けると、方法 13000 は再開され、判断ステップ 13010 で検査が行われて、架空の子ノード C_0 がコンテキスト・ノードとして選択されたかどうか判定される。そうである場合、実行は継続され、ステップ 13020 で、 C_0 がコンテキスト・ノードとして除外される。方法 13000 は、その後、ステップ 13015 で終了する。検査が 13010 で失敗した場合、この方法は、直接に終了ステップ 13015 へと進む。

10

【0136】

方法 13000 および 6000 は、そのキーワード検索での用途とは別に、階層的なデータの選択的な提示の手段としても使用することができる。すでに論じたように、実際の階層的なデータソースは、通常、ユーザが任意の所与の時点で見たいと希望するよりもずっと多くのデータを含んでいる。ユーザが、階層的なデータソースの閲覧を、そのデータ構造内部のノードを選択することによって行うとき、通常、提示アプリケーションは、選択したノードの下のサブツリーの中のデータ項目すべてを表示し、そのうちの一部は、しばしば、ユーザにとって興味のないものである可能性がある。その提示アプリケーションが、興味のないデータのふり分けを、ユーザの何らかのそれまでに観察された閲覧パターンに基づいて行えるのであれば、非常に望ましいはずである。説明した方法 13000 および 6000 は、この作業に非常に適している。 Y_i = ユーザが閲覧のために選択したルート・ノードと設定することにより、その方法によって識別されたコンテキスト・ノードの組は、ユーザにとって興味があり、表示することが好ましいようなノードを構成し、一方、コンテキスト・ノードとして識別されていない残りのノードは、好ましくはふり分けられる。

20

30

【0137】

親ノード P_j がルート・ノード Y_i からヒット・ノード X への有向パス上に存在する特殊な場合に関しては、コンテキスト・ノードとして含めるための Y_i の子ノード C_k の選択が、確率

【0138】

【数 26】

$$\text{Pr}[C_k | X \wedge P_j] = \frac{\text{Pr}[C_k \wedge X \wedge P_j]}{\text{Pr}[X \wedge P_j]} \quad \text{式 32}$$

40

【0139】

に基づくことを想起されたい。確率平均化によって、これらは、平均確率

【0140】

【数27】

$$\Pr_{mean}[C_k | X \wedge P_j] = \frac{\sum_h \Pr[C_{kh} \wedge X_h \wedge P_{jh}]}{\sum_h \Pr[X_h \wedge P_{jh}]}$$

$$\approx \frac{\sum_h \text{freq}(C_{kh}, P_{jh}, X_h)}{\sum_h \text{freq}(P_{jh}, X_h)}$$

式33

10

【0141】

によって置き換えられるが、ここで、 (P_{jh}, C_{kh}) は (P_j, C_k) に等しく、 (P_{jh}, X_h) は (P_j, X) に等しく、 P_{j0} 、 C_{k0} 、および X_0 （すなわち、 $h=0$ ）はそれぞれ、 P_j 、 C_k 、および X の別名とする。

【0142】

上の式が平均確率 $\Pr_{mean}[C_k | X \wedge P_j]$ の正確な近似であるためには、式33の右辺の分母が十分に大きくなければならない（たとえば、 $> f_{min}$ ）。そうでないときには、使用することのできる別の好ましい改善方法は、 $\Pr_{mean}[C_k | X \wedge P_j]$ の近似を、 $\Pr_{mean}[C_k | X' \wedge P_j]$ 、すなわち、 X ではなく X' に対する条件付き確率で行うことであり、ここで X' は Y_i から X への有向パス上に存在する X のすぐ上の親である。

20

【0143】

$\Pr_{mean}[C_k | X \wedge P_j]$ に対する近似の決定に使用するノード X' を識別するための方法22000のフローチャートを、図22に示している。方法22000は、ステップ22005で始まり、ここで、 X' はまず X に初期化される。次のステップ22010で、総和

【0144】

【数28】

$$\sum_h \text{freq}(P_{jh}, X'_h)$$

30

【0145】

が計算され、 D に割り当てられるが、ここでノード対 (P_{jh}, X'_h) は (P_j, X') に等しい。次いで、判断ステップ22015が続き、 D がある正の閾値定数 f_{min} 以上であるかどうかの検査が行われる。そうである場合は、方法22000は、ステップ22025で抜けて成功となる。判断ステップ22015が失敗した場合は、実行は別の判断ステップ22030へと進み、ここで検査が行われて、 X' が P_j のすぐ下の子であるかどうか判定される。そうである場合は、この方法はステップ22035で抜けて失敗となり、そうでない場合は、方法は継続して、ステップ22040で、 X' が、 P_j から X への有向パスに沿って存在するその親で置き換えられる。次いで、方法22000は、ループしてステップ22010へと戻る。

40

【0146】

方法22000が成功して、ノード X' および対応する値 D が得られた場合、 $\Pr_{mean}[C_k | X \wedge P_j]$ には、値

【0147】

【数29】

$$\Pr_{mean}[C_k | X \wedge P_j] \approx \frac{\sum_h \text{freq}(C_{kh}, P_{jh}, X'_h)}{D}$$

式34

50

【 0 1 4 8 】

が割り当てられる。

【 0 1 4 9 】

方法 2 2 0 0 0 を抜けて失敗になった場合には、 $\text{Pr}_{\text{mean}}[C_k | X \wedge P_j]$ への値の割り当ては、 C_k と Y_i の間の距離に基づいて

【 0 1 5 0 】

【 数 3 0 】

$$\text{Pr}_{\text{mean}}[C_k | X \wedge P_j] \approx \begin{cases} 1 & \text{if } \text{dist}(C_k, Y_i) \leq d_{\text{max}} \\ 0 & \text{otherwise} \end{cases} \quad \text{式 35}$$

10

【 0 1 5 1 】

または、 C_k と ヒット・ノード X の間の距離に基づいて

【 0 1 5 2 】

【 数 3 1 】

$$\text{Pr}_{\text{mean}}[C_k | X \wedge P_j] \approx \begin{cases} 1 & \text{if } \text{dist}(C_k, X) \leq d_{\text{max}} \\ 0 & \text{otherwise} \end{cases} \quad \text{式 36}$$

【 0 1 5 3 】

行われる。

20

【 0 1 5 4 】

$\text{Pr}_{\text{mean}}[C_k | X \wedge P_j]$ の近似が、最終的に式 3 4、式 3 5、または式 3 6 で行われたのに応じて、親ノード P_j に、コンテキストの子ノードが、 P_j から X への有向パス上に存在する子ノード C_1 以外ないことの、 P_j およびヒット・ノード X のもとの平均確率の計算は、それぞれ、式 3 7、式 3 8、または式 3 9 を用いて行われる。すなわち

【 0 1 5 5 】

【 数 3 2 】

$$\text{Pr}_{\text{mean}}[C_1 \text{ no sibling } | P_j \wedge X] \approx \frac{\sum_h \text{freq}(P_{j_h} \text{ has 1 child, } X'_h)}{D} \quad \text{式 37}$$

30

【 0 1 5 6 】

【 数 3 3 】

$$\text{Pr}_{\text{mean}}[C_1 \text{ no sibling } | P_j \wedge X] \approx \begin{cases} 0 & \text{if } \text{dist}(P, Y_i) + 1 \leq d_{\text{max}} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 38}$$

【 0 1 5 7 】

【 数 3 4 】

$$\text{Pr}_{\text{mean}}[C_1 \text{ no sibling } | P_j \wedge X] \approx \begin{cases} 0 & \text{if } \text{dist}(P, X) + 1 \leq d_{\text{max}} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 39}$$

40

【 0 1 5 8 】

ここで、 (P_{j_h}, X'_h) は (P_j, X') に等しく、 X' および D は方法 2 2 0 0 0 によって得られる。

【 0 1 5 9 】

親ノード P_j に対するコンテキストの子ノードの決定を、 P_j が Y_i からヒット・ノード X への有向パスに沿って存在する特別な場合について、確率平均化によって行うための好ましい手続きは、図 8 に示すものに非常に類似しており、これを図 1 4 に示している。

【 0 1 6 0 】

50

図14に示す方法14000は、ステップ14001で始まり、そこでは、概念上、架空の子ノード C_0 が作成され、実際の子ノード C_1, \dots, C_m からなるリストへと追加され、確率値 $Q_0 = Pr_{mean}[C_1 \text{ no sibling} | P_j X]$ が、式37、式38、または式39を用いて計算されて割り当てられ、次のステップ14005で、 C_1 を除く実際の子ノード C_k には、それに対応して、確率値 $Q_k = Pr_{mean}[C_k | X P_j]$ が、それぞれ、式34、式35、または式36を用いて計算されて割り当てられる。どの場合でも、ステップ14006がステップ14005に続き、方法7000をステップ7005で呼び出して(ステップ7001を飛ばして)、子ノード C_0, \dots, C_m のうちで1組のコンテキスト・ノードを選択する。方法7000を抜けると、方法14000は再開され、判断ステップ14010で検査が行われて、架空の子ノード C_0 がコンテキスト・ノードとして選択されたかが判定される。そうである場合、実行は継続され、ステップ14020で、 C_0 がコンテキスト・ノードとして除外される。方法14000は、その後、ステップ14015で終了する。検査が14010で失敗した場合、この方法は、直接に終了ステップ14015へと進む。

【0161】

前述の議論では、スキーマ・グラフ内に多くとも単一のヒット・ノードがある特別な場合においてコンテキスト・ノードを識別するための方法を説明している。これは、ユーザが単一の検索キーワードだけを入力するときの普通の想定状況である。そのキーワードがスキーマ・グラフ内の複数の場所に現れる場合には、これは1つより多くのヒットがあることを意味するが、ヒットそれぞれを別々に扱うことが好ましい。すなわち、説明した方法がスキーマ・グラフ内の第1のヒット・ノードに対して適用され、複数のコンテキストツリーがそのヒット・ノードに対して決定される。次いでその後で、同じ方法が残りのヒット・ノードのそれぞれに対して適用されて、新しい複数のコンテキストツリーが得られる、などである。ヒット・ノードすべてが処理されてしまうと、生成されたコンテキストツリーは、複数のヒット・ノードを包摂するとわかった場合には、スコアが付け直され、さらに重複しているコンテキストツリーは取り除かれる。次いで、残りのコンテキストツリーのリストは、(もしあれば)その新しいスコアに従って整列し直され、キーワード検索操作の結果としてユーザへと返される。

【0162】

しかし、ユーザが、ブール値のAND操作で結合された複数の検索キーワードを含む「すべてを見つける」(find all)キーワード検索操作を開始した場合は、キーワード・ヒットは、潜在的に、スキーマ・グラフ内の2つ以上のヒット・ノードに現れる可能性がある。コンテキストツリーを決定するためのより一般的な方法の説明を、ここで、そのような想定状況を扱うために行う。

【0163】

図11に、内部に複数のヒット・ノード11010、11020、および11025のある、スキーマ・グラフの例11000を示している。これらのヒット・ノードを X_1, \dots, X_n で表すことにする。当然、コンテキストツリーがヒット・ノードすべてを含むためには、Aで表している(ノード11005)、ヒット・ノードすべてを含む最も小さいサブツリーのルート・ノードが、コンテキスト・ノードとして、Aからヒット・ノードそれぞれへの有向パスに沿って存在するノードすべてとともに返されなければならない。したがって、ノード11015は、Aから X_2 へ(またAから X_3 へ)の有向パスに沿って存在するため、コンテキスト・ノードでなければならない。

【0164】

コンテキストツリー決定の方法の第1の、ボトム・アップのフェーズは、ノードAで始まり、走査は上向きに行われる。 Y_i をAまたはAの祖先とするが、ヒット・ノードのもとのその確率、すなわち $Pr[Y_i | X_1 \dots X_n]$ を評価することが、コンテキストツリーの可能なルート・ノードを決定するためには必要である。これを数学的に表現すると

【0165】

10

20

30

40

50

【数35】

$$\Pr[Y_i | X_1 \wedge \dots \wedge X_n] = \frac{\Pr[Y_i \wedge X_1 \wedge \dots \wedge X_n]}{\Pr[X_1 \wedge \dots \wedge X_n]} \quad \text{式 40}$$

【0166】

となる。

【0167】

この時点で、何らかの確率の独立性の仮定が必要であるが、理由は、式40の右辺の分子と分母はどちらも、 n の一般の値について既存の頻度テーブルから直接に求めまたは推定することが($n=2$ のときの分母を除いて)できないためである。妥当な1つの仮定は、 X_1 の組が、共通の祖先 Y_i のもとで、互いに独立というものである。言い換えれば、

$$\Pr[X_1 \wedge \dots \wedge X_n | Y_i] = \Pr[X_1 | Y_i] \cdot \dots \cdot \Pr[X_n | Y_i] \quad \text{式 41}$$

41

したがって、

【0168】

【数36】

$$\begin{aligned} \Pr[Y_i \wedge X_1 \wedge \dots \wedge X_n] &= \Pr[X_1 \wedge \dots \wedge X_n | Y_i] \Pr[Y_i] \\ &= \Pr[X_1 | Y_i] \cdot \dots \cdot \Pr[X_n | Y_i] \Pr[Y_i] \\ &= \frac{\Pr[X_1 \wedge Y_i] \cdot \dots \cdot \Pr[X_n \wedge Y_i]}{\Pr[Y_i]^{n-1}} \end{aligned} \quad \text{式 42}$$

【0169】

$\Pr[Y_i] = 0$ のときの特異性を取り除くために、 $\Pr[Y_i \wedge X_1 \wedge \dots \wedge X_n]$ は

【0170】

【数37】

$$\Pr[Y_i \wedge X_1 \wedge \dots \wedge X_n] = \begin{cases} 0 & \text{if } \Pr[Y_i] = 0 \\ \frac{\Pr[X_1 \wedge Y_i] \cdot \dots \cdot \Pr[X_n \wedge Y_i]}{\Pr[Y_i]^{n-1}} & \text{otherwise} \end{cases} \quad \text{式 43}$$

【0171】

同様に

【0172】

【数38】

$$\begin{aligned} \Pr[X_1 \wedge \dots \wedge X_n] &= \Pr[A \wedge X_1 \wedge \dots \wedge X_n] \\ &= \Pr[X_1 \wedge \dots \wedge X_n | A] \Pr[A] \\ &= \Pr[X_1 | A] \cdot \dots \cdot \Pr[X_n | A] \Pr[A] \\ &= \frac{\Pr[X_1 \wedge A] \cdot \dots \cdot \Pr[X_n \wedge A]}{\Pr[A]^{n-1}} \end{aligned} \quad \text{式 44}$$

【0173】

単一のヒット・ノードしかない場合のように、ヒット・ノードすべてのもとでの Y_i の生起確率の表現は、好ましくは、そのすぐ下の子ノードの確率によって漸進的に、確率平

10

20

30

40

50

均化が容易になるように、次のように行われる。すなわち、Wが、Y_iからAへの有向パスに沿うY_iのすぐ下の子ノードを表すとすると、

【0174】

【数39】

$$\Pr[Y_i | X_1 \wedge \dots \wedge X_n] =$$

$$\begin{cases} 1 & Y_i = A \\ 0 & Y_i \neq A, \Pr[W | X_1 \wedge \dots \wedge X_n] = 0 \\ \Pr_{mean}[Y_i | W \wedge X_1 \wedge \dots \wedge X_n] \cdot \Pr[W | X_1 \wedge \dots \wedge X_n] & otherwise \end{cases} \quad \text{式 45}$$

10

【0175】

ここで

【0176】

【数40】

$$\Pr_{mean}[Y_i | W \wedge X_1 \wedge \dots \wedge X_n] = \frac{\sum_h \Pr[Y_{ih} \wedge W_h \wedge X_{1h} \wedge \dots \wedge X_{nh}]}{\sum_h \Pr[W_h \wedge X_{1h} \wedge \dots \wedge X_{nh}]} \quad \text{式 46}$$

20

【0177】

ここで、対(Y_{ih}, W_h)は(Y_i, W)に等しく、l = 1, . . . , nに対して(W_h, X_{lh})は(W, X_l)に等しく、Y_{i0}、W₀、およびX_{l0} (h = 0)はそれぞれ、Y_i、W、およびX_lの別名とする。分子の総和の内側の項は、式43で置き換えることができる。分母の総和の内側の項も、WにY_iの役割をさせることにより、式43で置き換えることができ、こうして

【0178】

【数41】

$$\Pr_{mean}[Y_i | W \wedge X_1 \wedge \dots \wedge X_n] = \sum_h N_h / \sum_h D_h \quad \text{式 47}$$

30

【0179】

が得られ、ここで

【0180】

【数42】

$$N_h = \begin{cases} 0 & \text{if } \Pr[Y_{ih}] = 0 \\ \frac{\Pr[Y_{ih} \wedge X_{1h}] \cdots \Pr[Y_{ih} \wedge X_{nh}]}{\Pr[Y_{ih}]^{n-1}} & otherwise \end{cases}$$

40

$$\approx \begin{cases} 0 & \text{if } \text{freq}(Y_{ih}) = 0 \\ \frac{\text{freq}(Y_{ih}, X_{1h}) \cdots \text{freq}(Y_{ih}, X_{nh})}{\text{freq}(Y_{ih})^{n-1}} & otherwise \end{cases} \quad \text{式 48}$$

【0181】

【数 4 3】

$$D_h = \begin{cases} 0 & \text{if } \Pr[W_h] = 0 \\ \frac{\Pr[W_h \wedge X_{1h}] \cdots \Pr[W_h \wedge X_{nh}]}{\Pr[W_h]^{n-1}} & \text{otherwise} \end{cases}$$

$$\approx \begin{cases} 0 & \text{if } \text{freq}(W_h) = 0 \\ \frac{\text{freq}(W_h, X_{1h}) \cdots \text{freq}(W_h, X_{nh})}{\text{freq}(W_h)^{n-1}} & \text{otherwise} \end{cases}$$

式 49

10

【0182】

$\Pr_{\text{mean}}[Y_i | W, X_1, \dots, X_n]$ は、

【0183】

【数 4 4】

$$\sum_h D_h = 0.$$

20

【0184】

の場合には不定である。これが起きるとき、 $\Pr_{\text{mean}}[Y_i | W, X_1, \dots, X_n]$ への値の割り当ては、好ましくは、 Y_i からヒット・ノード X_1, \dots, X_n への距離に基づいて次のように行われる。すなわち

【0185】

【数 4 5】

$$\Pr_{\text{mean}}[Y_i | W \wedge X_1 \wedge \cdots \wedge X_n] = \begin{cases} 1 & \text{if } \min_{l=1, \dots, n} \text{dist}(Y_i, X_l) \leq d_{\text{max}} \\ 0 & \text{otherwise} \end{cases}$$

式 50

30

【0186】

ノード Y_i が、ヒット・ノードすべてを含むコンテキストツリーのルート・ノードである確率の計算を、 Y_i として選んだものすべてについて行うための方法 12000 のフローチャートを、図 12 に示している。方法 12000 は、ステップ 12001 で始まり、ここでは、ヒット・ノードすべてを含む、スキーマ・グラフ内の最小のサブツリーのルート・ノードが識別され、A で表される。次いで、実行は継続して、ステップ 12002 で、 Y_i が A に初期化され、したがって $\Pr'[Y_i | X_1, \dots, X_n] = 1$ となる。次のステップ 12005 で、式 45、式 47 を、式 48 および式 49、またはその代わりに式 50 とともに使用して、 $\Pr'[Z_j | X_1, \dots, X_n]$ の計算が、 Y_i の親ノード Z_j それぞれに対して行われる。ステップ 12005 の後、ステップ 12010 では、 $\Pr'[Y_i \text{ root} | X_1, \dots, X_n]$ の計算が、式

40

【0187】

【数 4 6】

$$\Pr'[Y_i \text{ root} | X_1 \wedge \cdots \wedge X_n] =$$

$$\Pr'[Y_i | X_1 \wedge \cdots \wedge X_n] - \sum_j \Pr'[Z_j | X_1 \wedge \cdots \wedge X_n]$$

式 51

【0188】

に従って行われる。

【0189】

50

次いで、方法12000はステップ12015へと進み、ここで、 Y_i の親ノード Z_j が選択される。ステップ12020に到達すると、方法12000は、ステップ12005で(ステップ12001および12002を飛ばして)再帰的に、ただし選択した親ノード Z_j が Y_i の役割を果たす状態で呼び出される。この再帰呼び出しが戻ると、実行が判断ステップ12025で再開され、ここで検査が行われて、 Y_i の親ノードすべてが処理されているかどうかの判定が行われる。そうである場合は、方法はステップ12030で終わり、そうでない場合は、方法は継続して、ステップ12015で Y_i の別の親ノード Z_j が、処理のために選択される。

【0190】

第2のトップ・ダウン走査のフェーズで、親ノード P_j がルート・ノード Y_i からどのヒット・ノードへの有向パスに沿っても存在しない一般の場合には、 P_j の子ノードがコンテキスト・ノードであるかどうかを判定するための方法は、これまでに単一のヒット・ノードしかない場合のために使用した方法13000とそのまま変わらない。

【0191】

親ノード P_j が Y_i から1つまたは複数のヒット・ノードへの有向パスに沿って存在する特別な場合には、単一のヒット・ノードの場合のために使用した方法を変更して、1つより多くの P_j の子ノードをコンテキスト・ノードとして含めなければならない可能性を見越しておく必要がある。1つのヒット・ノード X だけしか含まない場合には、子ノード C_k がコンテキスト・ノードであるかどうかの判定は、確率の値

$$\Pr[C_k | X, P_j]$$

に基づいて行われることを想起されたい。ここで、 X は P_j の子孫であるが、 C_k または C_k の子孫ではない。1つより多くのヒット・ノードがあるときの拡張は、子ノード C_k 、 $k = 1, \dots, m$ の選択プロセスを、親ノード P_j および P_j の子孫のであるヒット・ノード X_1 すべてのもとでの C_k の確率に基づくようにし、一方で P_j の子孫ではないヒット・ノードの影響を無視することである。当然、それ自身がヒット・ノードであり、または1つまたは複数のヒット・ノードの祖先である子ノード C_k すべては、コンテキスト・ノードでなければならない。一般性を失うことなく、これらの子ノードを、 C_1, \dots, C_r とし、ここで $1 \leq r \leq m$ とする。同様に、 C_1, \dots, C_r の子孫であるヒット・ノードの組を X_1, \dots, X_s とし、ここで $r \leq s \leq n$ とする。 $s = 1$ (したがって $r = 1$)の場合、この想定状況は、単一のヒット・ノードしかない場合に等しく、したがって、この場合に関して説明した方法14000を使用することができる。 $s > 1$ の場合に関して一般化するための好ましい実装形態で採用する方法は、項 $\Pr[C_k | X, P_j]$ を、式

【0192】

【数47】

$$\sum_{l=1}^s \Pr[C_k | X_l \wedge P_j]$$

【0193】

で置き換えることであり、これは、確率平均化の後、

【0194】

【数48】

$$Q_k = \sum_{l=1}^s \Pr_{\text{mean}}[C_k | X_l \wedge P_j]$$

式52

【0195】

となり、ここで $\Pr_{\text{mean}}[C_k | X, P_j]$ は、式33で定義した通りである。 Q_k は、任意の X_1 、 $1 = 1, \dots, s$ に対して $\text{freq}_{\text{mean}}(P_j, X_1) = 0$ の場合には、不定となり、ここで

【0196】

10

20

30

40

50

【数 49】

$$freq_{mean}(P_j, X_l) = \sum_h freq(P_{jh}, X_{lh})$$

【0197】

であり、 (P_{jh}, X_{lh}) は (P_j, X_l) に等しい。 $freq_{mean}(P_j, X_l)$ が非零ではあるが小さな数である（たとえば、 $< f_{min}$ ）のときでも、 Q_k の推定は、頻度テーブルからは十分な精度で行うことができない。単一のヒット・ノードを含む場合のように、この問題の克服は、ヒット・ノード X_1, \dots, X_s を、そのヒット・ノードの一部またはすべてをその祖先で置き換えた新しいノードの組 S で置き換え、その後で Q_k を S の要素によって定義し直すことによって行うことができる。

10

【0198】

図 21 のフローチャートで示している方法 21000 が、好ましくは、ヒット・ノード X_1, \dots, X_s を置き換えるこのノードの新しい組 S を決定するのに使用される。方法 21000 は、ステップ 21005 で始まり、ここで、ヒット・ノードの初期の組 X_1, \dots, X_s は S で表される。次のステップ 21010 で、 S のうちの処理済みでない要素 X_p が選択される。次いで、判断ステップ 21015 が続き、ここでは検査が行われて、 $freq_{mean}(P_j, X_p)$ が、ある閾値定数 f_{min} 以上であるかが判定される。そうである場合、 X_p は組 S の中に保持され、方法 21000 は判断ステップ 21020 へと続き、ここで検査が行われて S の中の要素すべてが処理されているかが判定される。1 つまたは複数の処理済みでない要素が残っている場合は、実行はステップ 21010 へと戻って、 S の別の要素が処理のために選択される。他方、要素すべてが処理されている場合は、方法はステップ 21025 で終わって成功になる。

20

【0199】

ここで判断ステップ 21015 へと戻る。検査の条件が失敗した場合、別の判断ステップ 21030 が続き、これにより、選択したノード X_p が P_j の子ノードであるかが検査される。そうである場合は、方法 21000 はステップ 21040 で終わって失敗となり、そうでない場合は、ステップ 21035 が続く。ステップ 21035 で、 S の中の要素 X_p は、 P_j から X_p への有向パスに沿って存在するその親 X'_p で置き換えられ、 X'_p の子孫すべてが S から取り除かれる。次いで、実行はステップ 21020 へと進む。

30

【0200】

上で説明した方法 21000 が成功で戻った場合は、結果となる組 S の中の要素を使用して、値 Q_k の計算が子ノード C_1, \dots, C_r それぞれに対して行われる。すなわち

【0201】

【数 50】

$$Q_k = \sum_{X \in S} Pr_{mean}[C_k | X \wedge P_j]$$

式 53

【0202】

しかし、方法 21000 が失敗で戻った場合は、値 Q_k の決定は、好ましくは、 C_k とルート・ノード Y_i の間の距離から、すなわち

40

【0203】

【数 51】

$$Q_k = \begin{cases} 1 & \text{dist}(C_k, Y_i) \leq d_{max} \\ 0 & \text{otherwise} \end{cases}$$

式 54

【0204】

またはこの代わりに、 C_k とヒット・ノード X_1, \dots, X_n の間の距離から、すなわち

【0205】

50

【数52】

$$Q_k = \begin{cases} 1 & \min_{i=1, \dots, n} \text{dist}(C_k, X_i) \leq d_{\max} \\ 0 & \text{otherwise} \end{cases} \quad \text{式 55}$$

【0206】

のように行われる。

【0207】

また、単一のヒット・ノード X がある場合、親ノード P_j およびヒット・ノード X のもとで、 P_j の1つの子ノード C_1 だけが生起する確率

$$\text{Pr}[C_1 \text{ no sibling} | X \in P_j]$$

の評価が必要であることを想起されたい。ここで、 C_1 は X であり、または X の祖先である。

【0208】

この量を現在の想定状況へと一般化する際、2つの可能性、すなわち、 $r = 1$ である特別な場合、および $r > 1$ であるより一般的な場合が生じる可能性がある。前者の例を図15に示しているが、そこには、ノード15005をルートとするサブツリーの内部にある3つのヒット・ノード15030、15035、および15040がある。しかし、3つのヒット・ノードはすべて、ノード15005の単一の子ノード15010の下にある。

【0209】

この特別な場合を扱うための1つのアプローチは、項 $\text{Pr}[C_1 \text{ no sibling} | X \in P_j]$ の式

【0210】

【数53】

$$Q_0 = \sum_{l=1}^s \text{Pr}_{\text{mean}}[C_1 \text{ no sibling} | X_l \in P_j] \quad \text{式 56}$$

【0211】

による置き換えを、量 Q_k を式52で使用したのと類似のしかたで行うことである。上の Q_k の場合と同様に、式56は、任意の X_l 、 $l = 1, \dots, s$ に対して $\text{freq}_{\text{mean}}(P_j, X_l) = 0$ の場合には、不定となる可能性がある。したがって、 Q_0 に割り当てられる実際の値は、方法21000が成功で戻った場合は、方法21000から得られた組 S に基づく。すなわち

【0212】

【数54】

$$Q_0 = \sum_{X \in S} \text{Pr}_{\text{mean}}[C_1 \text{ no sibling} | X \in P_j] \quad \text{式 57}$$

【0213】

そうでなく方法21000が失敗であった場合は、 Q_0 への値の割り当ては、 P_j のルート・ノード Y_i からの距離に基づいて

【0214】

【数55】

$$Q_0 = \begin{cases} 0 & \text{dist}(P_j, Y_i) + 1 \leq d_{\max} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 58}$$

【0215】

またはこの代わりに、 P_j とヒット・ノード X_1, \dots, X_n の間の距離に基づいて

【0216】

10

20

30

40

【数56】

$$Q_0 = \begin{cases} 0 & \min_{i=1, \dots, n} \text{dist}(P_j, X_i) + 1 \leq d_{\max} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 59}$$

【0217】

行われる。

【0218】

親ノード P_j の子ノード C_k の組のうちでコンテキスト・ノードの識別を $r = 1, s > 1$ の場合に関して行うための、図16のフローチャートに示す方法16000は、単一の
 ヒット・ノードの場合のための方法14000に非常に似ている。方法16000は、ス
 テップ16001で始まり、そこでは、概念上、架空の子ノード C_0 が作成され、実際
 の子ノード C_1, \dots, C_m からなるリストへと追加され、式57、式58、または式5
 9の中で定義された値 Q_0 が割り当てられ、次のステップ16005で、 C_1 を除く実際
 の子ノード C_k には、それに対応して、それぞれ、式53、式54、または式55の中
 で定義された値 Q_k が割り当てられる。ステップ16006はステップ16005に続き、
 方法7000をステップ7005で呼び出して（ステップ7001を飛ばして）、子ノ
 ード C_0, \dots, C_m のうちで1組のコンテキスト・ノードを選択する。方法7000を
 抜けると、方法16000は再開され、判断ステップ16010で検査が行われて、架空
 の子ノード C_0 がコンテキスト・ノードとして選択されたかが判定される。そう
 である場合、実行は継続され、ステップ16020で、 C_0 がコンテキスト・ノードとして
 除外される。方法16000は、その後、ステップ16015で終了する。検査が160
 10で失敗した場合、方法16000は、直接に終了ステップ16015へと進む。

【0219】

$r > 1$ （したがって $s > 1$ ）である一般の場合では、 $r = 1$ の場合に使用した $\text{Pr}[C_1 \text{ no sibling} | X_1 \dots X_s P_j]$ に類似した量は

【0220】

【数57】

$$\sum_{Y \in S} \text{Pr}[C_1 \wedge \dots \wedge C_r \wedge \neg C_{r+1} \wedge \dots \wedge \neg C_m | X \wedge P_j] \quad 30$$

【0221】

であり、ここで、 S は、方法21000が成功で抜けた場合に、これによって返される
 組である。残念ながら、総和の中にある確率の推定は、既存の頻度テーブルからは容易
 に行うことができない。したがって、わずかに異なる式がその代わりに使用される。子ノ
 ード $C_k, 1 \leq k \leq r$ それぞれにルートがあるサブツリーの中にない組 S の要素を $1 \leq l \leq s_k$
 に対する H_{kl} で表すとし、ここで $s_k = |S|$ である。子ノード $C_k, 1 \leq k \leq r$
 それぞれに対し、次のものが計算される。すなわち

【0222】

【数58】

$$Q_k = \sum_{l=1}^{s_k} \text{Pr}_{\text{mean}}[C_k | H_{kl} \wedge P_j] \quad \text{式 60}$$

【0223】

上の式の量の背後にある理論的根拠は、 $C_k, 1 \leq k \leq r$ すべてにわたって足し合わせ
 ると、子ノード C_1, \dots, C_r がともに生起する確率を近似する量が得られるという
 ものである（が、 $s_k > 1$ である値をとることがあるため、真の確率ではない）。 $r = 1$ の場
 合でのように、方法21000が失敗で戻った場合、 Q_k は、 P_j のルート・ノード Y_i
 からの距離から、 $1 \leq k \leq r$ に対して

【0224】

10

20

30

40

50

【数59】

$$Q_i = \begin{cases} 0 & \text{dist}(P_j, Y_i) + 1 \leq d_{\max} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 61}$$

【0225】

またはこの代わりに、 P_j とヒット・ノード X_1, \dots, X_n の間の距離から

【0226】

【数60】

$$Q_k = \begin{cases} 0 & \min_{i=1, \dots, n} \text{dist}(P_j, X_i) + 1 \leq d_{\max} \\ 1 & \text{otherwise} \end{cases} \quad \text{式 62} \quad 10$$

【0227】

のように得られる。

【0228】

子ノードの組 C_1, \dots, C_m のうちでコンテキスト・ノードを選択するための方法 17000 を、ここで、 $r > 1$ 、 $s > 1$ の一般の場合に関して、図 17 のフローチャートを参照して説明する。方法 17000 は、ステップ 17001 で始まり、ここで、子ノード C_k 、 $1 \leq k \leq r$ それぞれには、式 60、式 61、または式 62 を用いて計算された値 Q_k が割り当てられる。ステップ 17005 がこれに続き、ここでは、残りの子ノードには、それに対応して、式 53、式 54、または式 55 を用いて計算された値 Q_k が、それぞれ、割り当てられる。次のステップ 17010 で、値 Q_k が、子ノードすべてにわたって総計され、 T と表される。方法 17000 は継続して、ステップ 17015 で、そのサブツリーにヒット・ノードを含む子ノードすべて、すなわち、 C_k 、 $1 \leq k \leq r$ がコンテキスト・ノードとして選択される。次のステップ 17020 で、割り当てられた値が残りの子ノードのうちで最高のノード C_k もコンテキスト・ノードとして選択される。値が最大で同じ子ノードが 1 つより多く存在する場合は、そのようなノードはすべて、コンテキスト・ノードとして選択される。次いで、それまでにコンテキスト・ノードとして選択された子ノードすべての割り当てられた値の総和が、ステップ 17025 で計算され、 S と表される。次いで、実行は判断ステップ 17030 へと進み、このポイントで、子ノード C_k すべてがコンテキスト・ノードとして選択されている場合、方法 17000 は、ステップ 17040 で終了する。しかし、コンテキスト・ノードとしてまだ選択されていない 1 つまたは複数の子ノード C_k がある場合は、方法 17000 は、別の判断ステップ 17035 へと続く。ステップ 17035 で、検査が行われて、 $S < T/2$ であるかどうか確かめられ、そうである場合、方法 17000 はここでもステップ 17040 で終了する。 $S < T/2$ である場合は、実行はステップ 17020 へと戻り、ここでノードのコンテキスト・ノードとしての選択がさらに行われる。

【0229】

以上の説明では、階層的なデータ構造内のキーワード検索を行うときに遭遇する異なる段階および動作の想定状況を扱うための様々な方法を提示している。こうした方法は、単一の全体的な手続き 18000 へと組み込まれるが、これは、図 2 のステップ 2010 を詳細にしたものであり、図 19 および図 20 にそれぞれ示す下位の手続き 19000 および 20000 を含む図 18 のフローチャートで示している。方法 18000 は、判断ステップ 18005 で始まり、ここで、検査が行われて、スキーマ・グラフ内に複数のヒット・ノードがあるかどうかの判定が行われる。そうである場合、実行はステップ 18015 へと進み、ここで方法 20000 が呼び出され、そうでない場合は、実行はステップ 18010 へと進み、ここで方法 19000 が呼び出される。どちらの場合でも、方法 20000 または 19000 は、コンテキストツリーからなるリストとともに戻り、それぞれには、スコアが関連づけられている。以下は、方法 19000 の詳細な説明であり、これに方法 20000 のそれが続く。

【0230】

方法19000は、ステップ19001で始まり、ここで、方法10000が呼び出されて、ヒット・ノードXの祖先ノードである可能なルート・ノード Y_i からなるリストが決定される。 Y_i それぞれは、可能なコンテキストツリーのルート・ノードである。また、方法10000では、値 $S_i = Pr' [Y_i | X]$ の計算がノード Y_i それぞれに対して行われる。次いで、方法19000は継続して、ステップ19005で、その前のステップで決定されたノード Y_i が処理のために選択される。次のステップ19010で、方法19000内部にある下位プロセスである方法38000が呼び出されて、ノード Y_i をルートとするサブツリーの中にあるコンテキスト・ノードが識別される。次いで、方法19000は継続して、ステップ19030で、識別されたコンテキスト・ノードすべてを含み、 Y_i をルート・ノードとするコンテキストツリーが構築される。このツリーには、ステップ19001で計算された S_i のスコアが割り当てられる。次いで、方法19000は、判断ステップ19035へと進む。ステップ19001で得られたノード Y_i すべてが処理されている場合は、この方法はステップ19040で終わり、そうでない場合は、この方法はステップ19005へと戻って、別のノード Y_i の処理が行われる。

10

【0231】

方法19000の内部で呼び出される方法38000は、ステップ38010で始まり、ここで、ノード Y_i が、まず P_j に割り当てられる。実行は、判断ステップ38015へ、次いで、 P_j が Y_i からヒット・ノードXへの有向パス上に存在しない場合は、38020へと進む。ステップ38020で、方法13000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードが選択される。続くステップ38025では、方法38000の再帰的な呼び出しが、ステップ38020で(ステップ38010および38015を飛ばして)、コンテキスト・ノードとして選択された葉でない子ノード C_k それぞれに対して行われ、 C_k に P_j の役割をさせることによって、さらに別のコンテキスト・ノードがその子孫のうちで識別される。そのような子ノードすべてに対する呼び出しが戻ると、方法38000は、ステップ38040で終了する。また、方法38000は、 P_j に子ノードがない場合、またはその葉でない子ノードのうちのどれかが、ステップ38020でコンテキスト・ノードとして選択されていない場合には、直接に終了ステップ38040へと進む。

20

【0232】

判断ステップ38015は、 P_j が Y_i からXへの有向パス上に存在する場合に成功し、その場合、実行はステップ38045へと進む。ここで、方法14000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードが選択されるが、 C_1 は、 P_j からXへの有向パス上に存在する子ノードを表す。続くステップ38050で、方法38000の再帰的な呼び出しが、ステップ38015で(ステップ38010を飛ばして)、コンテキスト・ノードとして選択された葉でない子ノード C_k それぞれに対して行われ、 C_k に P_j の役割をさせることによって、さらに別のコンテキスト・ノードがその子孫のうちで識別される。そのような子ノードすべてに対する呼び出しが戻ると、方法38000は、ステップ38040で終了する。

30

【0233】

方法20000は、ステップ20001で始まり、ここで、方法12000が呼び出されて、ヒット・ノード X_1, \dots, X_n の祖先ノードである可能なルート・ノード Y_i からなるリストが決定される。 Y_i それぞれは、可能なコンテキストツリーのルート・ノードである。また、方法12000では、値 $S_i = Pr' [Y_i | X_1 \dots X_n]$ の計算がノード Y_i それぞれに対して行われる。次いで、方法20000は継続して、ステップ20005で、その前のステップで決定されたノード Y_i が処理のために選択される。次のステップ20010で、方法20000内部にある下位プロセスである方法39000が呼び出されて、ノード Y_i をルートとするサブツリーの中にあるコンテキスト・ノードが識別される。次いで、方法20000は継続して、ステップ20060で、識別されたコンテキスト・ノードすべてを含み、 Y_i をルート・ノードとするコンテキストツ

40

50

リーが構築される。このツリーには、ステップ20001で計算された S_i のスコアが割り当てられる。次いで、方法20000は、判断ステップ20065へと進む。ステップ20001で得られたノード Y_i すべてが処理されている場合は、この方法はステップ20070で終わり、そうでない場合は、この方法はステップ20005へと戻って、別のノード Y_i の処理が行われる。

【0234】

方法19000の内部で呼び出される方法39000は、ステップ39010で始まり、ここで、ノード Y_i が、まず P_j に割り当てられる。実行は、判断ステップ39015へ、次いで、 P_j のところのサブツリーのルートにおいてヒット・ノードがない場合は、ステップ39020へと進む。ステップ39020で、方法13000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードが選択される。続くステップ39025では、方法39000の再帰的な呼び出しが、ステップ39020で(ステップ39010および39015を飛ばして)、コンテキスト・ノードとして選択された葉でない子ノード C_k それぞれに対して行われ、 C_k に P_j の役割をさせることによって、さらに別のコンテキスト・ノードがその子孫のうちで識別される。そのような子ノードすべてに対する呼び出しが戻ると、方法39000は、ステップ39060で終了する。また、方法39000は、 P_j に子ノードがない場合、またはその葉でない子ノードのうちどれかが、ステップ39020でコンテキスト・ノードとして選択されていない場合には、直接に終了ステップ39060へと進む。

【0235】

判断ステップ39015は、 P_j をルートとするサブツリーの内部に1つまたは複数のヒット・ノードがある場合に成功し、その場合、実行は別のステップ39030へと進む。 P_j の下のサブツリーの中に単一のヒット・ノードしかない場合は、この判断ステップは失敗し、実行はステップ39035へと進み、そうでない場合は、実行はまた別の判断ステップ39040へと進む。判断ステップ39040で、検査が行われて、 P_j の下のヒット・ノードすべてが、その子ノードのうちの下にしかないかが判定される。そうである場合は、実行はステップ39045へと進み、そうでない場合は、実行はステップ39050へと進む。ステップ39050では、 C_1, \dots, C_r で1つまたは複数のヒット・ノードが下にある P_j の子ノードを表すとして、方法17000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードが選択される。しかし、判断ステップ39040からステップ39045へと導かれた場合は、方法16000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードを選択することが、 C_1 をそのサブツリーの中にヒット・ノードを含む、 P_j の唯一の子ノードとして行われる。

【0236】

ここで39035へと戻り、 P_j からその唯一の子孫のヒット・ノードへのパスが、その子ノード C_1 を通過するとする。方法14000が呼び出されて、 P_j の子ノードのうちで1組のコンテキスト・ノードが選択される。

【0237】

ステップ39035、39045、および39050のそれぞれが完了すると、方法39000のそれ自体の再帰的な呼び出しが、ステップ39015で(ステップ39010を飛ばして)、コンテキスト・ノードとして選択された葉でない子ノード C_k それぞれに対して行われ、 C_k に P_j の役割をさせることによって、さらに別のコンテキスト・ノードがその子孫のうちで識別される。そのような子ノードすべてに対する呼び出しが戻ると、方法39000は、ステップ39060で終了する。

【0238】

説明のための例

ここで、好ましい実装形態の動作の説明を、下の階層的なXMLデータソースの例によって行う。このXMLソースは、「XYZ」という名称の企業に関する、Webアドレス、支社の名称および場所、各支社での販売製品の範囲などのデータを含んでいる。このX

10

20

30

40

50

MLデータのスキーマ・グラフ表現を、図23に示している。

【0239】

XMLソース

```

<company> <! - - 企業 - - >
  <name>XYZ</name> <! - - 名称 - - >
  <web>http://www.xyz.com</web> <! - - We
bアドレス - - >
  <description> <! - - 説明 - - >
    Company founded in 1999 specialising
    in hi-tech consumers electronics <! - - 19
99年設立のハイテク家電製品に特化した企業 - - >
  </description>
  <branch> <! - - 支社 - - >
    <name>North Ryde</name> <! - - 名称 - - >
    <phone>0291230000</phone> <! - - 電話番号 - -
>
    <address> <! - - 住所 - - >
      <number>1</number> <! - - 番地 - - >
      <street>Lane Cove</street> <! - - 街区 C
ove通り - - >
      <city>Sydney</city> <! - - 都市 シドニー - - >
      <country>Australia</country> <! - - 国
オーストラリア - - >
    </address>
    <manager> <! - - 支社長 - - >
      <firstName>Jim</firstName> <! - - 名 - - >
      <lastName>Smith</lastName> <! - - 姓 - - >
      <email>jsmith@xyz.com</email> <!@電子
メール - - >
    </manager>
    <product> <! - - 製品 - - >
      <id>1</id> <! - - ID - - >
      <name>Plasma TV</name> <! - - 名称プラズマTV
- - >
      <price>$10000</price> <! - - 価格 - - >
      <supplier>JEC</supplier> <! - - メーカー - -
>
      <stock>10</stock> <! - - 在庫数 - - >
    </product>
    <product> <! - - 製品 - - >
      <id>2</id> <! - - ID - - >
      <name>Mp3 player</name> <! - - 名称 MP3プ
レイヤー - - >
      <price>$500 </price> <! - - 価格 - - >
      <supplier>HG</supplier> <! - - メーカー - - >
      <stock>20</stock> <! - - 在庫数 - - >
    </product>
  </branch>
  <branch> <! - - 支社 - - >
    <name>Morley</name> <! - - 名称 - - >

```

```

    < phone > 0 8 9 1 2 3 0 0 0 0 < / phone >  < ! - - 電話番号 - -
>
    < address >  < ! - - 住所 - - >
      < number > 1 < / number >  < ! - - 番地 - - >
      < street > Russel < / street >  < ! - - 街区  Russ
e l - - >
      < city > Perth < / city >  < ! - - 都市 パース - - >
      < country > Australia < / country >  < ! - - 国
オーストラリア - - >
    < / address >
    < manager >  < ! - - 支社長 - - >
      < firstName > Ted < / firstName >  < ! - - 名 - - >
      < lastName > White < / lastName >  < ! - - 姓 - - >
      < email > twhite@xyz.com < / email >  < ! @電子
メール - - >
    < / manager >
    < product >  < ! - - 製品 - - >
      < id > 3 < / id >  < ! - - ID - - >
      < name > Video phone < / name >  < ! - - 名称 ビデオ
電話 - - >
      < price > $ 2 0 0 0 < / price >  < ! - - 価格 - - >
      < supplier > NVC < supplier >  < ! - - メーカー - - >
      < stock > 15 < / stock >  < ! - - 在庫数 - - >
    < / product >
    < product >  < ! - - 製品 - - >
      < id > 4 < / id >  < ! - - ID - - >
      < name > PDA < / name >  < ! - - 名称 PDA - - >
      < price > $ 1 0 0 0 < / price >  < ! - - 価格 - - >
      < supplier > LP < / supplier >  < ! - - メーカー - - >
      < stock > 50 < / stock >  < ! - - 在庫数 - - >
    < / product >
  < / branch >
< / company >

```

図 23 で、各ノードの次に示す整数は、そのノードに割り当てられている一意の ID である。このデータソースの既存のビューが 3 つあるとする。第 1 のものは、この企業の名称、説明、および Web アドレスを表示しているビューである。第 2 のものは、この企業の支社およびその場所のリストであり、最後に第 3 のビューには、各支社での製品の範囲 (line) がリストされている。これらのビューのスキーマ・グラフ表現を、図 24、図 25、および図 26 にそれぞれ示してある。これらのビューの結果として、生起 27000、共起 28000、葉共起 29000、および単独子共起 30000 の頻度テーブルを、図 27、図 28、図 29、および図 30 にそれぞれ示してある。同時生起頻度テーブルは、3 次元であるが、5 つの別々の 2 次元テーブル 31000、32000、33000、34000、および 35000 によって示してある。図 31 は、 $P_j =$ ノード 1 であるテーブル内のエントリ $f_{req}(C_k, P_j, X)$ を含んでいる。同様に、図 32、図 33、図 34、および図 35 それぞれは、 $P_j =$ ノード 3、ノード 8、ノード 9、およびノード 10 であるエントリをそれぞれ含む。ここに示す頻度テーブルすべてにおいて、図 28 に見られるような、項目 28005 などの空のセルは、それに関連する頻度を保存する必要のない無効なノードの組合せを表す。

【 0 2 4 0 】

ユーザが、自分のいる市の中で特定の製品を見つけることを希望しているとする。ユー

10

20

30

40

50

ずは、その製品の名前、「MP3プレイヤー」および市の名前「シドニー」を入力し、両方の名前に対するキーワード検索を行う。図23からわかるように、これから、2つのヒット・ノード $X_1 = \text{ノード19}$ および $X_2 = \text{ノード13}$ が得られる。このキーワード検索操作に対する可能なコンテキストツリーを決定するために、システム4000では、図18の方法18000を呼び出す。ヒット・ノードが1つより多くあることから、方法18000は、その後、方法20000の呼び出しをステップ18015で行う。方法20000は、今度は、方法12000の呼び出しをステップ20001で行って、結果となるコンテキストツリーのルート・ノードの役をするノード Y_i のリストが得られる。

【0241】

方法12000では、まず、ステップ12001で、ノード3を、ヒット・ノード X_1 と X_2 をどちらも含む最小のサブツリーのルート・ノードとして識別する。したがって、 $A = \text{ノード3}$ である。次いで、方法12000は、再帰手続きを始めて、このヒット・ノードのもとで、 A およびその祖先のそれぞれに対する生起確率の値を計算する。ノード A では

$$\Pr' [A | X_1, X_2] = 1$$

である。 $Y_1 = \text{ノード1}$ 、すなわちノード A の親では、式47、式48、および式49を用いて

【0242】

【数61】

$$\Pr_{\text{mean}} [\text{node 1} | A \wedge X_1 \wedge X_2] =$$

$$\frac{\text{freq}(\text{node 1}, \text{node 19}) \text{freq}(\text{node 1}, \text{node 13}) \text{freq}(\text{node 3})}{\text{freq}(\text{node 3}, \text{node 19}) \text{freq}(A, \text{node 13}) \text{freq}(\text{node 1})}$$

【0243】

したがって

$$\Pr' [\text{node 1} | X_1, X_2] = 0$$

であり

$$\Pr' [A \text{ root} | X_1, X_2] = 1$$

である。

【0244】

このため、方法12000は、ノード A をコンテキストツリーの単一のルート・ノード候補として抜ける。このコンテキストツリーには1というスコアが割り当てられる。方法12000の完了後、方法20000では、第2の、トップ・ダウン走査のフェーズが続けられ、そこでは、ルート・ノード $Y_1 = A$ の子孫が処理されて、それらのうちでコンテキスト・ノードが識別される。このフェーズはステップ20010で始まり、ここで、 P_j がまずノード3になるように設定される。このノードは、2つの異なる子ノードの下にある、ヒット・ノード X_1 および X_2 の祖先であることから、実行は、結局、方法39000のステップ39050へと進み、ここで方法17000が呼び出されて、コンテキスト・ノードがその子のうちで決定される。方法17000の結果、ノード3の子ノード6~10それぞれに割り当てられる値 Q_1, \dots, Q_5 は、次の通りである。

【0245】

【数 6 2】

$$\begin{aligned}
 Q_1 &= \Pr_{\text{mean}}[\text{node 6} | X_1 \wedge P_j] + \Pr_{\text{mean}}[\text{node 6} | X_2 \wedge P_j] \\
 &= \frac{\text{freq}(\text{node 6, node 3, node 19})}{\text{freq}(\text{node 3, node 19})} + \frac{\text{freq}(\text{node 6, node 3, node 13})}{\text{freq}(\text{node 3, node 13})} \\
 &= \frac{1}{1} + \frac{1}{1} \\
 &= 2
 \end{aligned}$$

10

$$\begin{aligned}
 Q_2 &= \Pr_{\text{mean}}[\text{node 7} | X_1 \wedge P_j] + \Pr_{\text{mean}}[\text{node 7} | X_2 \wedge P_j] \\
 &= \frac{\text{freq}(\text{node 7, node 3, node 19})}{\text{freq}(\text{node 3, node 19})} + \frac{\text{freq}(\text{node 7, node 3, node 13})}{\text{freq}(\text{node 3, node 13})} \\
 &= 2
 \end{aligned}$$

$$Q_3 = \Pr_{\text{mean}}[\text{node 8} | X_1 \wedge P_j]$$

20

$$\begin{aligned}
 &= \frac{\text{freq}(\text{node 8, node 3, node 19})}{\text{freq}(\text{node 3, node 19})} \\
 &= 0
 \end{aligned}$$

$$Q_4 = \Pr_{\text{mean}}[\text{node 9} | X_1 \wedge P_j] + \Pr_{\text{mean}}[\text{node 9} | X_2 \wedge P_j]$$

$$\begin{aligned}
 &= \frac{\text{freq}(\text{node 9, node 3, node 19})}{\text{freq}(\text{node 3, node 19})} + \frac{\text{freq}(\text{node 9, node 3, node 13})}{\text{freq}(\text{node 3, node 13})} \\
 &= 1
 \end{aligned}$$

30

$$Q_5 = \Pr_{\text{mean}}[\text{node 10} | X_2 \wedge P_j]$$

$$\begin{aligned}
 &= \frac{\text{freq}(\text{node 10, node 3, node 13})}{\text{freq}(\text{node 3, node 13})} \\
 &= 0
 \end{aligned}$$

40

【0 2 4 6】

したがって、降順でソートされた組 Q_1, \dots, Q_5 は、 $\{Q_1, Q_2, Q_4, Q_3, Q_5\}$ であり、その総計は $T = 5$ となる。方法 17000 によって選択されたコンテキスト・ノードの組は、こうして、ノード 6、ノード 7 ($Q_1 + Q_2 > T/2$ であるため)、およびノード 8、ノード 10 (ヒット・ノードの祖先であるため) を含む。次いで、方法 39000 は、ステップ 39055 で再開して、それ自体を再帰的に呼び出して、子のある選択されたノードそれぞれに対してコンテキストの子ノードの識別を行う。

【0 2 4 7】

P_j = ノード 8 の場合、実行は、ノード 8 に単一の子孫のヒット・ノード (ノード 13

50

)があることから、ステップ39035へと進み、このポイントで方法14000が呼び出されて、コンテキスト・ノードが子ノードの組11~14のうちで識別される。方法14000の結果、 P_j の子ノード11、12、および14にそれぞれ割り当てられる確率の値 Q_1 、 Q_2 、 Q_4 は次の通りである。

【0248】

【数63】

$$Q_1 = \Pr_{\text{mean}}[\text{node 11} | X_2 \wedge P_j]$$

$$= \frac{\text{freq}(\text{node 11}, \text{node 8}, \text{node 13})}{\text{freq}(\text{node 8}, \text{node 13})} \quad 10$$

=1

$$Q_2 = \Pr_{\text{mean}}[\text{node 12} | X_2 \wedge P_j]$$

$$= \frac{\text{freq}(\text{node 12}, \text{node 8}, \text{node 13})}{\text{freq}(\text{node 8}, \text{node 13})} \quad 20$$

=1

$$Q_4 = \Pr_{\text{mean}}[\text{node 14} | X_2 \wedge P_j]$$

$$= \frac{\text{freq}(\text{node 14}, \text{node 8}, \text{node 13})}{\text{freq}(\text{node 8}, \text{node 13})}$$

=1

【0249】 30

さらに、方法14000では、架空の子ノード C_0 に対する値 Q_0 の計算も行われる。

【0250】

【数64】

$$Q_0 = \Pr_{\text{mean}}[\text{node 13 no sibling} | X_2 \wedge P_j]$$

$$= \frac{\text{freq}(\text{node 8 has 1 child}, \text{node 13})}{\text{freq}(\text{node 8}, \text{node 13})}$$

=0 40

【0251】

したがって、降順でソートされた確率の値の組は、 $\{Q_1, Q_2, Q_4, Q_0\}$ であり、その総計は $T=3$ となる。方法14000によって選択されたコンテキスト・ノードの組は、こうして、ノード11、ノード12、ノード14 ($Q_1 + Q_2 + Q_4 > T/2$ かつ $Q_1 = Q_2 = Q_4$ であるため)、およびノード13 (ヒット・ノードの祖先であるため)を含む。

【0252】

これに似た実行パスを、 $P_j =$ ノード10の場合にはたどることになり、似た結果が得られる。ノード10のコンテキストの子ノードの組は、ノード18~22である。したが 50

って、このコンテキストツリーのスキーマ・グラフ3600は、図36に示すものようになり、ヒット・ノード19および13、ならびにコンテキスト・ノード3、6～8、10～14、18～22を含む。これらのノードの表現するデータ項目を含む、ユーザへと返される実際のコンテキストツリーは、次のようなものである。

【0253】

```

    <branch> <! - - 支社 - - >
      <name>North Ryde</name> <! - - 名称 - - >
      <phone>0291230000</phone> <! - - 電話番号 - - >
    >
      <address> <! - - 住所 - - > 10
      <number>1</number> <! - - 番地 - - >
      <street>Lane Cove</street> <! - - 街区 Co
v e 通り - - >
      <city>Sydney</city> <! - - 都市 シドニー - - >
      <country>Australia</country> <! - - 国 オ
ー ストラリア - - >
      </address>
      <product> <! - - 製品 - - >
      <id>1</id> <! - - ID - - >
      <name>Plasma TV</name> <! - - 名称 プラズマTV 20
- - >
      <price>$10000</price> <! - - 価格 - - >
      <supplier>JEC</supplier> <! - - メーカー - - >
      <stock>10</stock> <! - - 在庫数 - - >
      </product>
      <product> <! - - 製品 - - >
      <id>2</id> <! - - ID - - >
      <name>Mp3 player</name> <! - - 名称 MP3プレ
イヤー - - >
      <price>$500</price> <! - - 価格 - - > 30
      <supplier>HG</supplier> <! - - メーカー - - >
      <stock>20</stock> <! - - 在庫数 - - >
      </product>
    </branch>
    <branch> <! - - 支社 - - >
      <name>Morley</name> <! - - 名称 - - >
      <phone>0891230000</phone> <! - - 電話番号 - - >
    >
      <address> <! - - 住所 - - >
      <number>1</number> <! - - 番地 - - > 40
      <street>Russel</street> <! - - 街区 Russe
l - - >
      <city>Perth</city> <! - - 都市 パース - - >
      <country>Australia</country> <! - - 国 オ
ー ストラリア - - >
      </address>
      <product> <! - - 製品 - - >
      <id>3</id> <! - - ID - - >
      <name>Video phone</name> <! - - 名称 ビデオ電
話 - - > 50

```

```

< price > $ 2 0 0 0 < / price > < ! - - 価格 - - >
< supplier > NVC < supplier > < ! - - メーカー - - >
< stock > 15 < / stock > < ! - - 在庫数 - - >
< / product >
< product > < ! - - 製品 - - >
< id > 4 < / id > < ! - - ID - - >
< name > PDA < / name > < ! - - 名称 PDA - - >
< price > $ 1 0 0 0 < / price > < ! - - 価格 - - >
< supplier > LP < / supplier > < ! - - メーカー - - >
< stock > 50 < / stock > < ! - - 在庫数 - - >
< / product >
< / branch >

```

10

産業上の利用可能性

上記から、説明した配置がコンピュータおよびデータ処理の産業に対して、特に複数の検索からの情報を提示することに関して適用可能であることは明らかである。

【0254】

上記では、本発明の一部の実施形態だけを説明しており、これへの変更および/または改変は、本発明の範囲および趣旨から逸脱することなく行うことができ、これら実施形態は例示的であり、制限的ではない。

【0255】

20

(オーストラリア限定)この明細書の文脈では、「含む」(comprising)が意味するのは、「主として、しかし必ずしももっぱらにではなく含む」(including principally but not necessarily solely)または「有する」(having)または「含む」(including)であって、「だけから構成される」(consisting only of)ではない。「comprise」や「comprises」など、「comprising」という語の変種の意味は、これに対応して変わる。

【図面の簡単な説明】

【0256】

【図1】スキーマ・グラフの例の図である。

30

【図2】キーワード検索方法のフローチャートである。

【図3A】スキーマ・グラフ内の親ノードの2つの例を示す図である。

【図3B】スキーマ・グラフ内の親ノードの2つの例を示す図である。

【図4】サーバおよびクライアントのコンピュータからなるネットワークの図である。

【図5】ルート・ノードからヒット・ノードへの有向パスに沿って存在しない親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図6】ルート・ノードからヒット・ノードへの有向パスに沿って存在しない親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための別の方法のフローチャートである。

40

【図7】ルート・ノードからヒット・ノードへの有向パスに沿って存在する親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図8】ルート・ノードからヒット・ノードへの有向パスに沿って存在する親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための別の方法のフローチャートである。

【図9】2つの同じサブツリーをスキーマ・グラフの例の図である。

【図10】確率平均化によるコンテキスト・ノード識別方法の第1の、ボトム・アップ横断フェーズのフローチャートである。

【図11】複数のヒット・ノードのあるスキーマ・グラフの例の図である。

50

【図12】確率平均化による、複数のヒット・ノードを含むコンテキスト・ノード識別方法の第1の、ボトム・アップ横断フェーズのフローチャートである。

【図13】確率平均化による、ルート・ノードからヒット・ノードへの有向パスに沿って存在しない親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図14】確率平均化による、ルート・ノードからヒット・ノードへの有向パスに沿って存在する親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図15】子孫のヒット・ノードがすべて単一の子ノードの下に位置する親ノードの例の図である。

10

【図16】複数のヒット・ノードすべてが単一の子ノードの下に位置する場合に対する、確率平均化による、ルート・ノードからヒット・ノードへの有向パスに沿って存在する親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図17】複数のヒット・ノードが複数の子ノードの下に位置する場合に対する、確率平均化による、ルート・ノードからヒット・ノードへの有向パスに沿って存在する親ノードの1組の子ノードのうちでコンテキスト・ノードを識別するための方法のフローチャートである。

【図18】1つまたは複数のヒット・ノードが存在する可能性のある、コンテキスト・ノードを識別するための方法のフローチャートである。

20

【図19】単一のヒット・ノードを含む場合に対するコンテキストツリーを構築するための方法のフローチャートである。

【図20】複数のヒット・ノードを含む場合に対するコンテキストツリーを構築するための方法のフローチャートである。

【図21】観測頻度がヒット・ノードからなる元の組での観測頻度よりも高いヒット・ノードからなる代替の組を構築するための方法のフローチャートである。

【図22】観測頻度がヒット・ノードからなる組よりも高いそのヒット・ノードからなる組の祖先を選択するための方法のフローチャートである。

【図23】スキーマ・グラフの例の図である。

【図24】データ・ビューの例のスキーマ・グラフの図である。

30

【図25】データ・ビューの別の例のスキーマ・グラフの図である。

【図26】データ・ビューのさらに別の例のスキーマ・グラフの図である。

【図27】図24、25、および26のデータ・ビューから生じる生起頻度テーブルの図である。

【図28】図24、25、および26のデータ・ビューから生じる共起頻度テーブルの図である。

【図29】図24、25、および26のデータ・ビューから生じる葉の共起頻度テーブルの図である。

【図30】図24、25、および26のデータ・ビューから生じる単独子の共起頻度テーブルの図である。

40

【図31】図24、25、および26のデータ・ビューから生じる同時生起頻度テーブルのある部分の図である。

【図32】図24、25、および26のデータ・ビューから生じる同時生起頻度テーブルの別の部分の図である。

【図33】図24、25、および26のデータ・ビューから生じる同時生起頻度テーブルのさらに別の部分の図である。

【図34】図24、25、および26のデータ・ビューから生じる同時生起頻度テーブルのさらに別の部分の図である。

【図35】図24、25、および26のデータ・ビューから生じる同時生起頻度テーブルのさらに別の部分の図である。

50

【図36】2つのキーワードを含むキーワード検索操作の結果として返されるコンテキストツリーのスキーマ・グラフの図である。

【図37】記載した配置を実施できる汎用コンピュータの図式的な構成図である。

【図38】図19に示す単一のヒット・ノードを含む場合に対するコンテキストツリーを構築するための方法内部の下位プロセスのフローチャートである。

【図39】図20に示す複数のヒット・ノードを含む場合に対するコンテキストツリーを構築するための方法内部の下位プロセスのフローチャートである。

【図1】

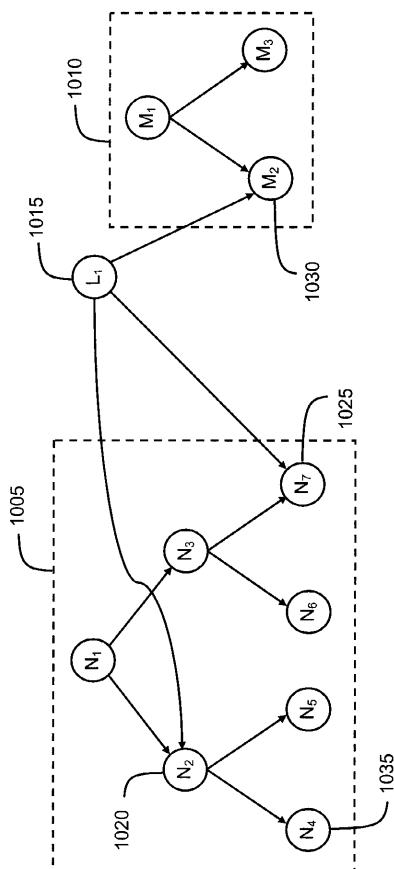


Fig. 1

【図2】

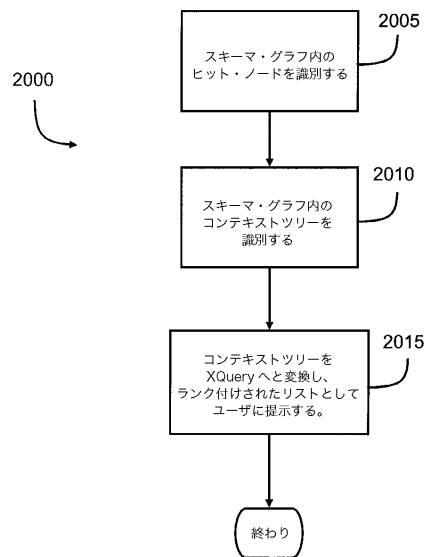


Fig. 2

【図3A】

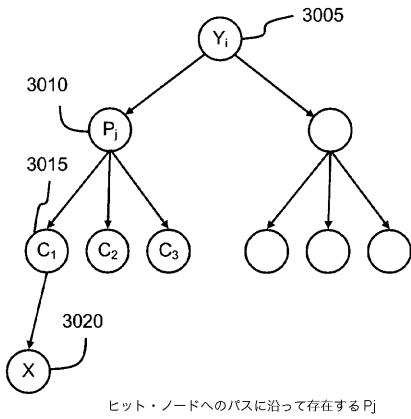


Fig. 3A

【図3B】

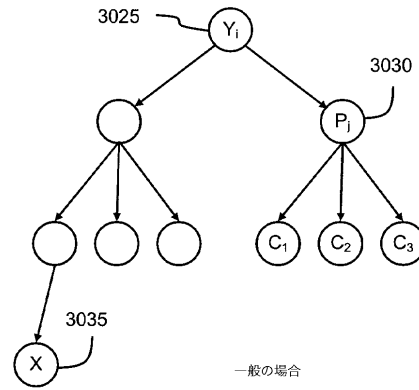


Fig. 3B

【図4】

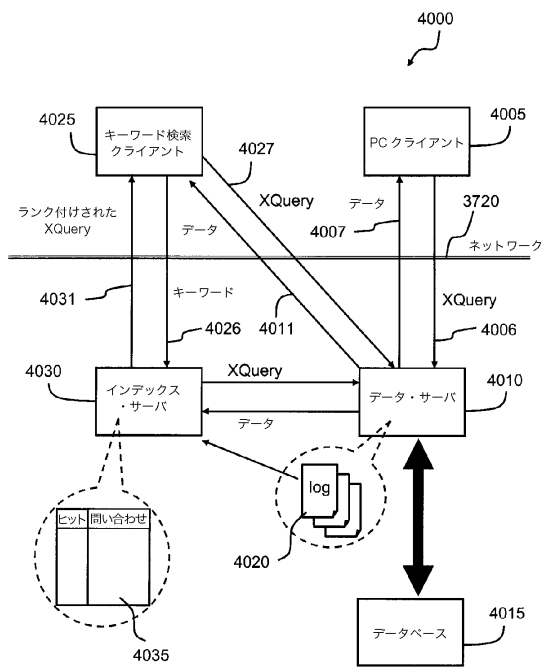


Fig. 4

【図5】

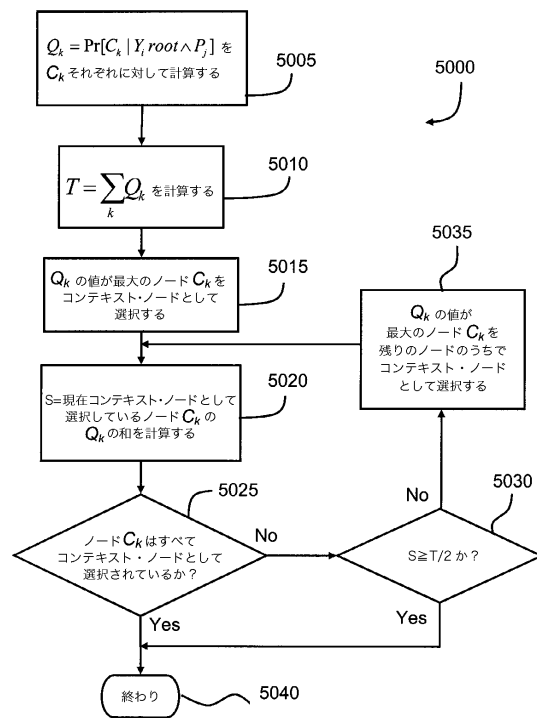


Fig. 5

【 図 6 】

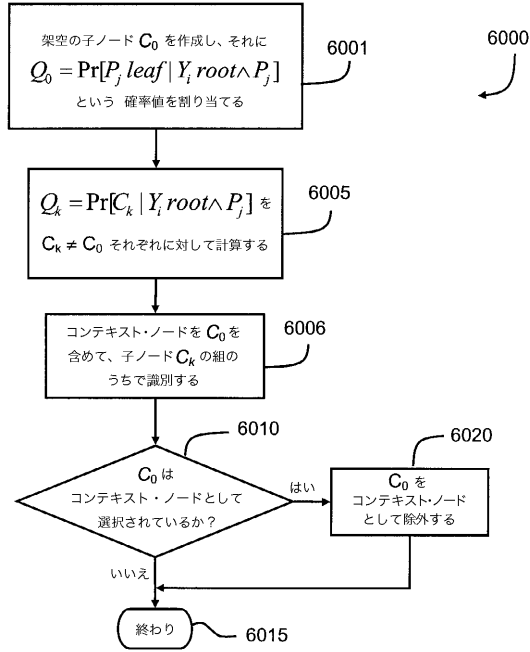


Fig. 6

【 図 7 】

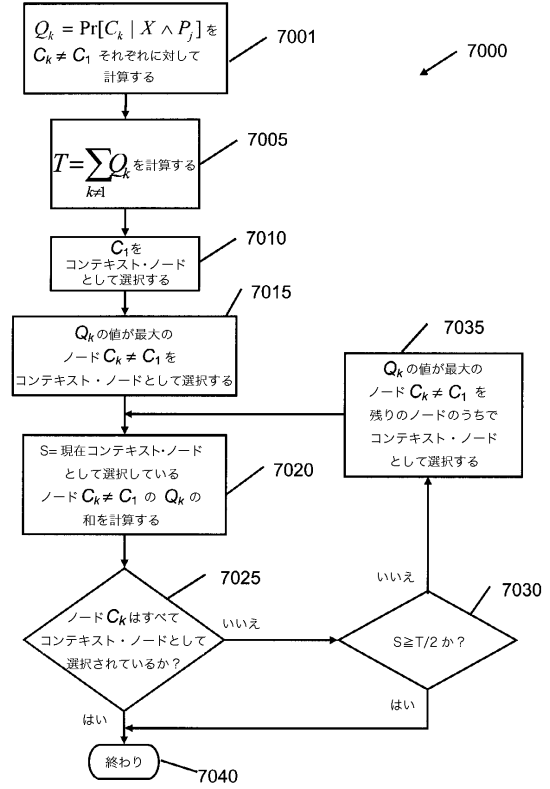


Fig. 7

【 図 8 】

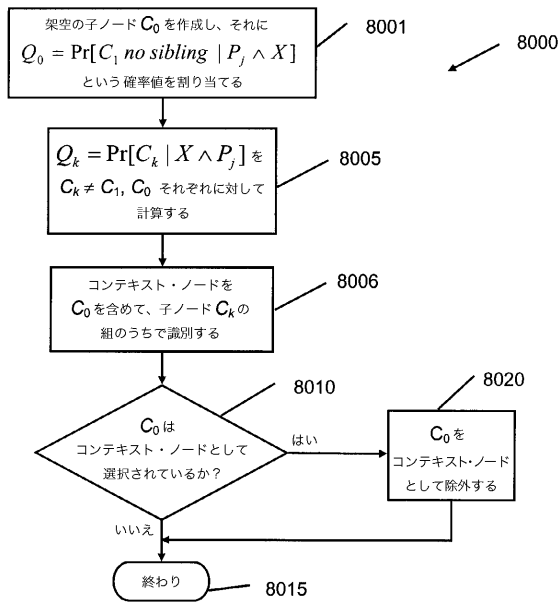


Fig. 8

【 図 9 】

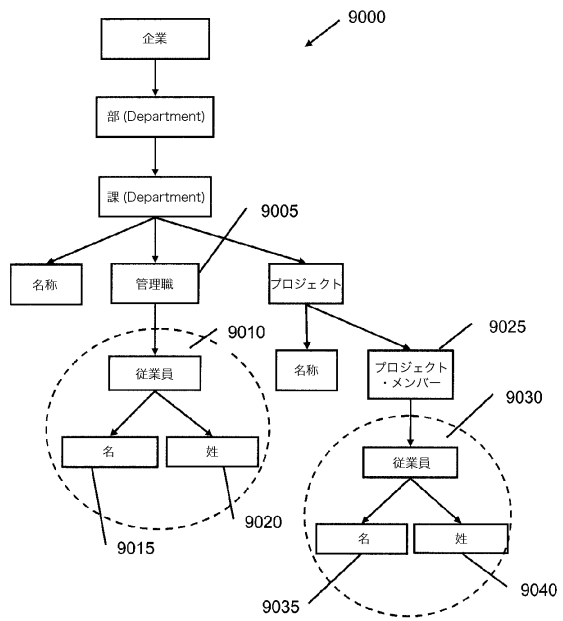


Fig. 9

【 図 1 0 】

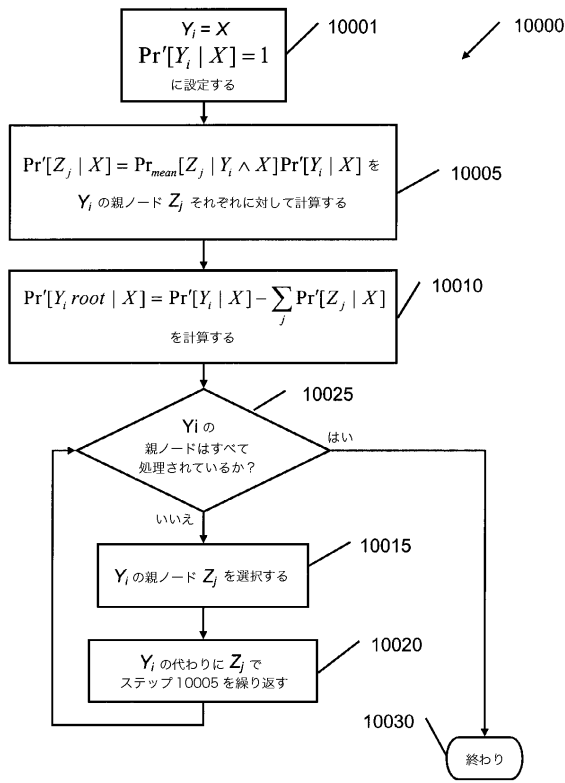


Fig. 10

【 図 1 1 】

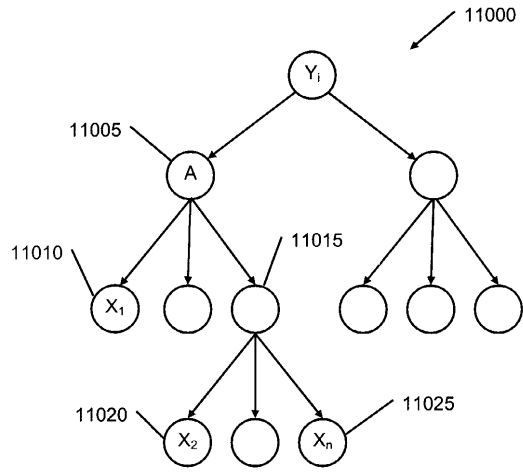


Fig. 11

【 図 1 2 】

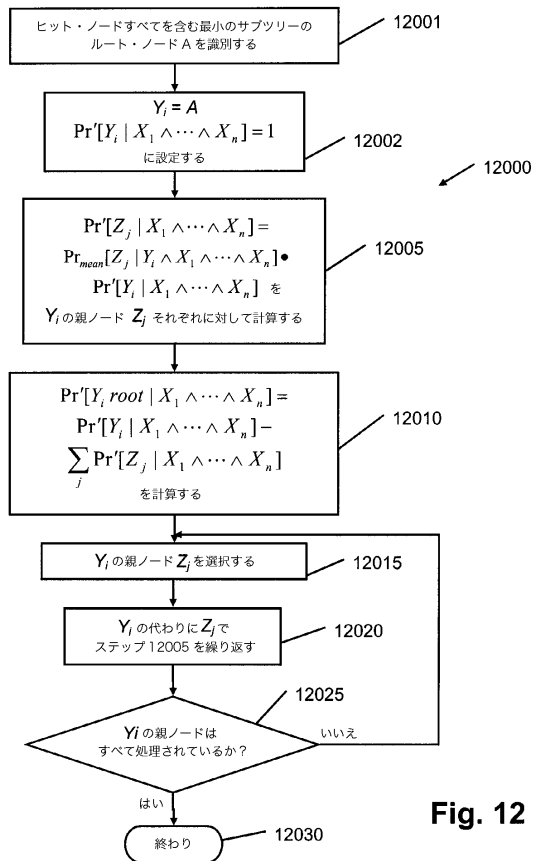


Fig. 12

【 図 1 3 】

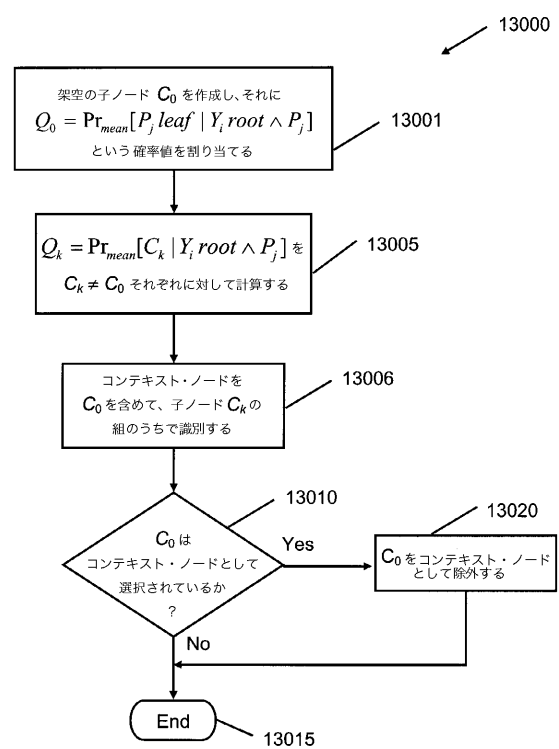


Fig. 13

【 図 1 4 】

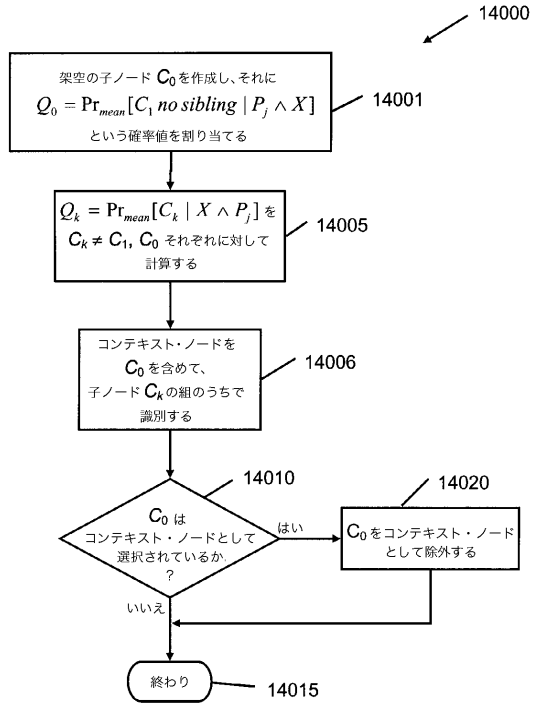


Fig. 14

【 図 1 5 】

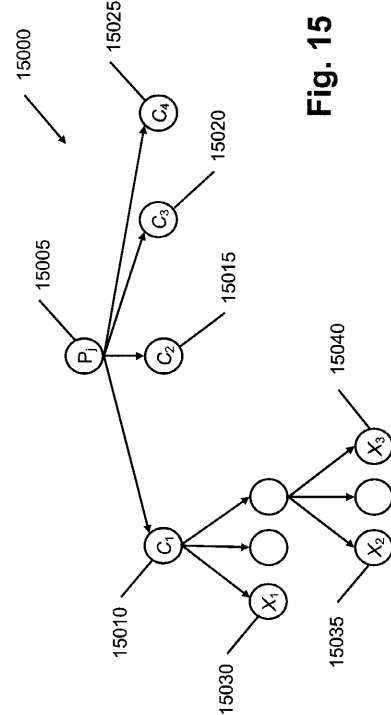


Fig. 15

【 図 1 6 】

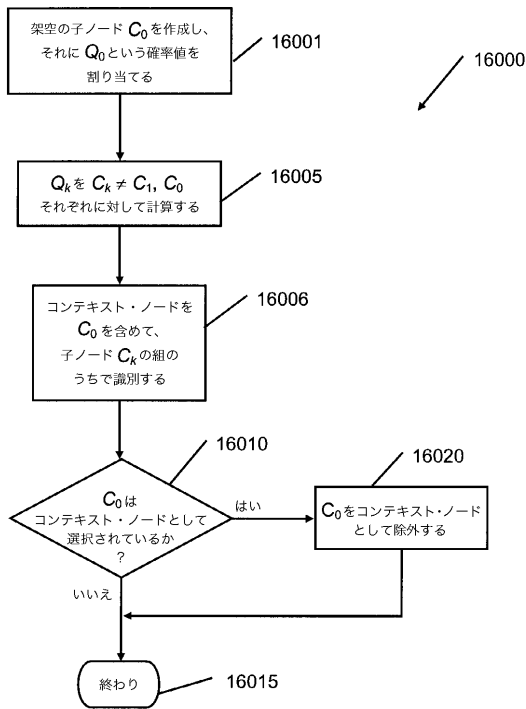


Fig. 16

【 図 1 7 】

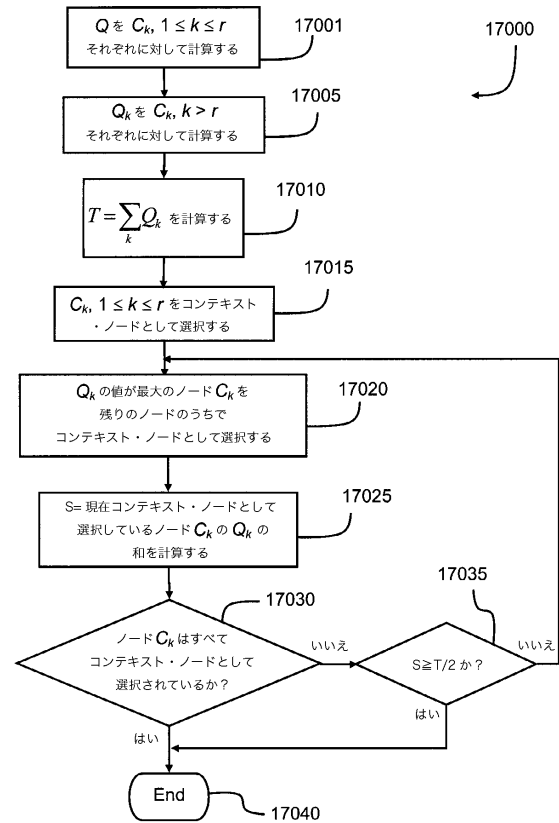


Fig. 17

【図18】

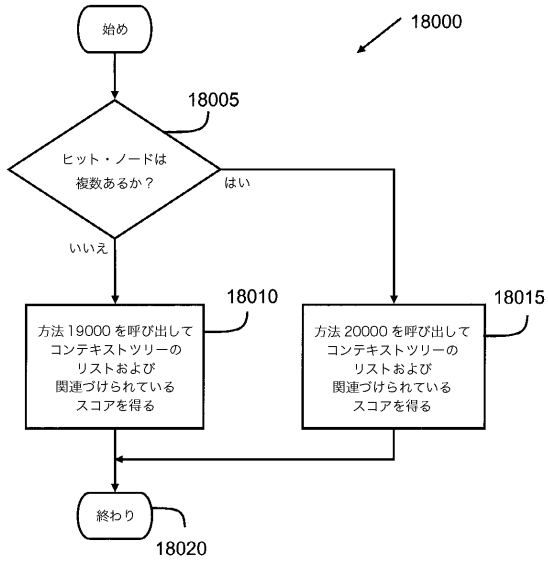


Fig. 18

【図19】

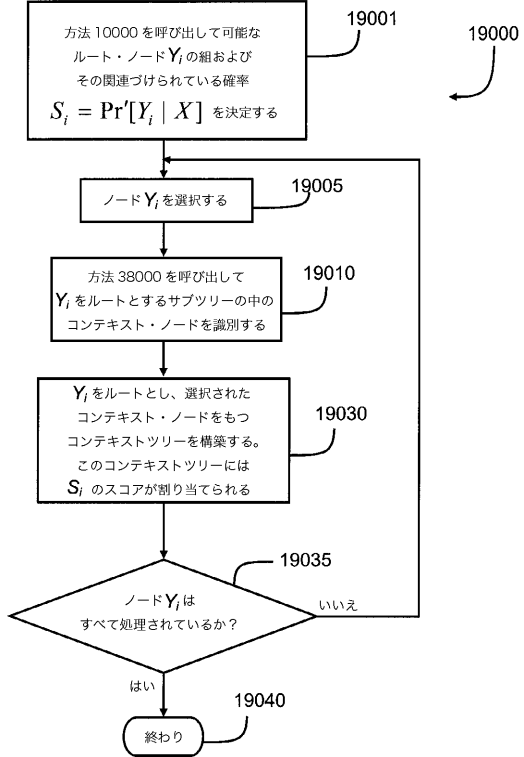


Fig. 19

【図20】

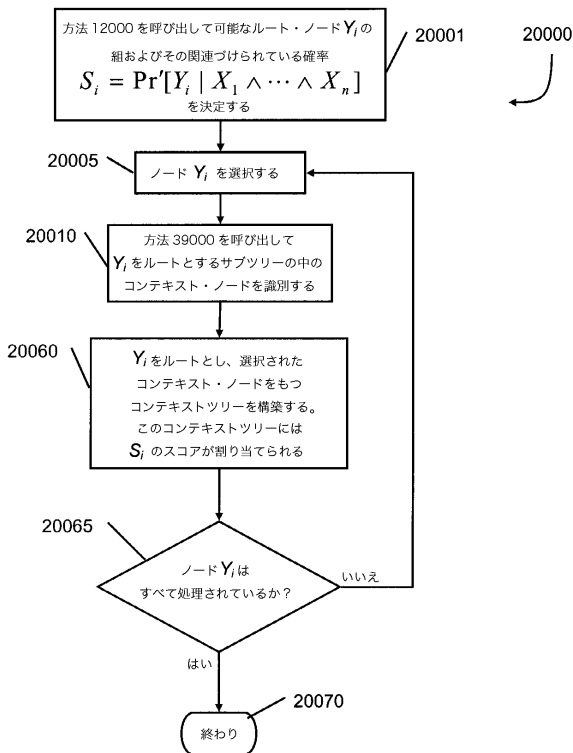


Fig. 20

【図21】

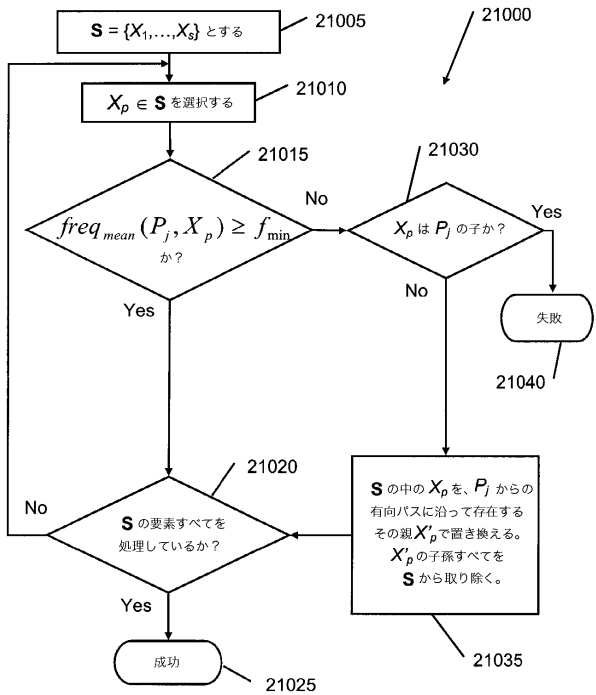


Fig. 21

【 図 2 2 】

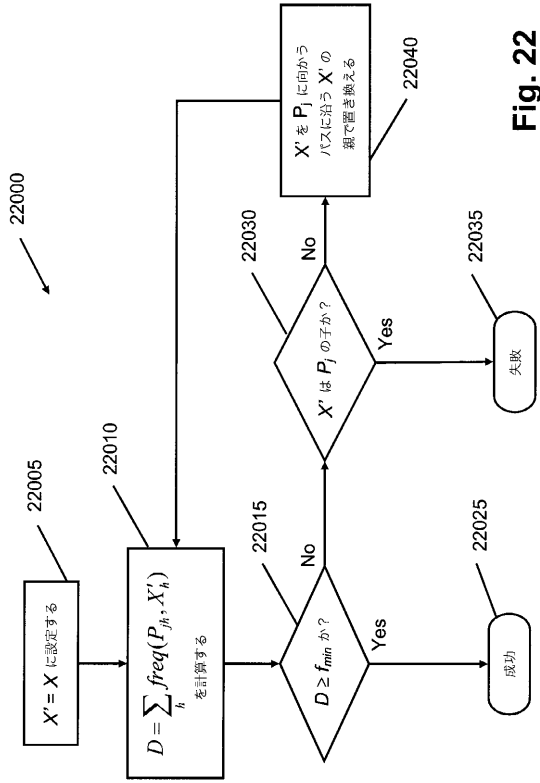


Fig. 22

【 図 2 3 】

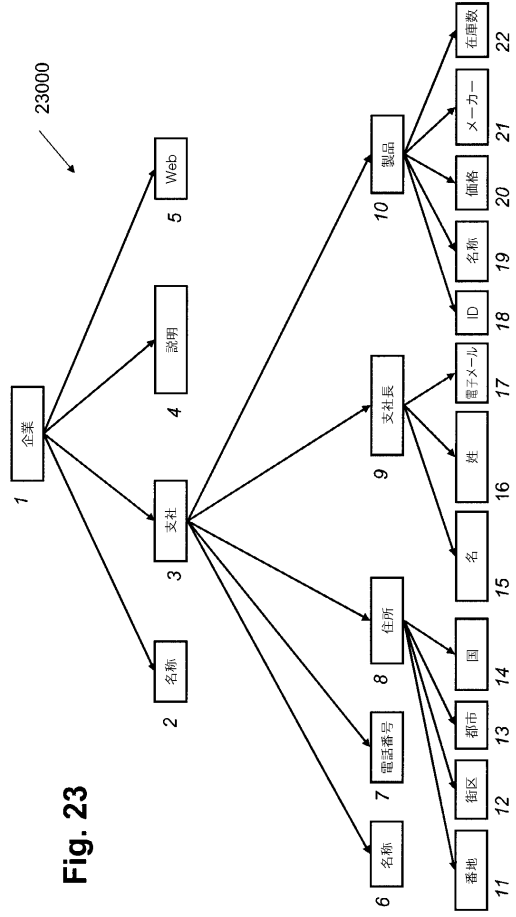


Fig. 23

【 図 2 4 】

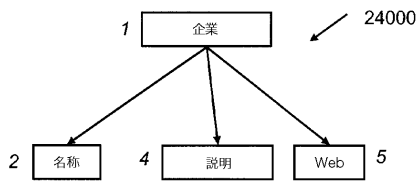


Fig. 24

【 図 2 5 】

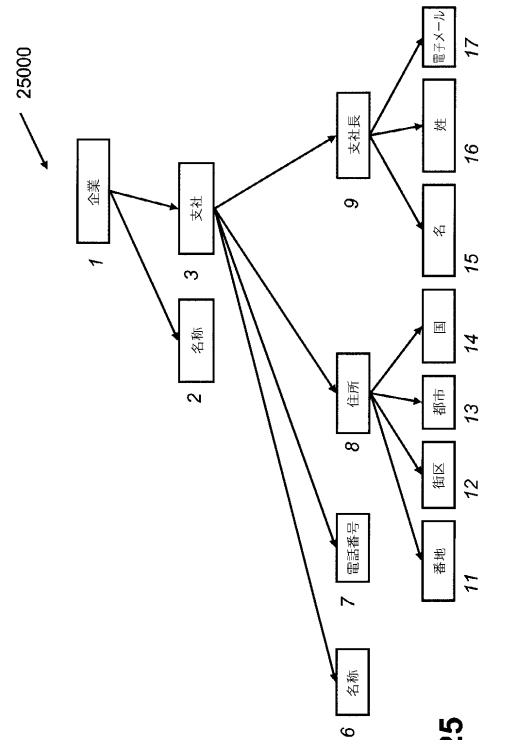


Fig. 25

【 図 26 】

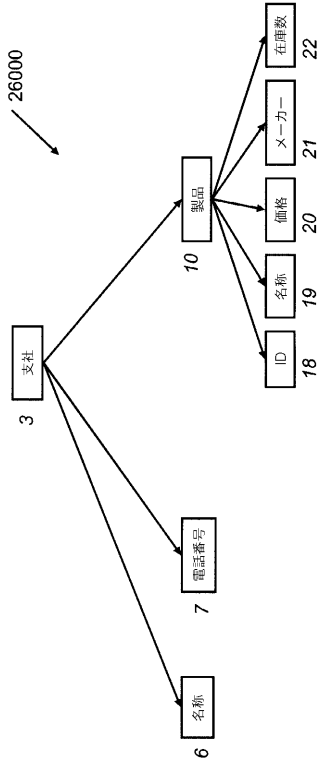


Fig. 26

【 図 27 】

X	$freq(X)$
1	2
2	2
3	2
4	1
5	1
6	2
7	2
8	1
9	1
10	1
11	1
12	1
13	1
14	1
15	1
16	1
17	1
18	1
19	1
20	1
21	1
22	1

Fig. 27

【 図 28 】

$freq(Y_i, X)$

X	Y_i				
	1	3	8	9	10
2	2				
3	1				
4	1				
5	1				
6	1	2			
7	1	2			
8	1	1			
9	1	1			
10	0	1			
11	1	1	1		
12	1	1	1		
13	1	1	1		
14	1	1	1		
15	1	1		1	
16	1	1		1	
17	1	1		1	
18	0	1			1
19	0	1			1
20	0	1			1
21	0	1			1
22	0	1			1

Fig. 28

【 図 30 】

$freq(P_j \text{ has 1 child}, X)$

X	P_j				
	1	3	8	9	10
2	0				
3	0				
4	0				
5	0				
6	0	0			
7	0	0			
8	0	0			
9	0	0			
10	0	0			
11	0	0	0		
12	0	0	0		
13	0	0	0		
14	0	0	0		
15	0	0		0	
16	0	0		0	
17	0	0		0	
18	0	0			0
19	0	0			0
20	0	0			0
21	0	0			0
22	0	0			0

Fig. 30

【 図 29 】

$freq(Y_i, P_j \text{ leaf})$

P_j	Y_i	
	1	3
3	0	
8	0	0
9	0	0
10	0	0

Fig. 29

【 図 3 1 】

$freq(C_k, P_j, X): P_j = \text{ノード}1$

X	C _k			
	2	3	4	5
2		1	1	1
3	1		0	0
4	1	0		1
5	1	0	1	
6	1		0	0
7	1		0	0
8	1		0	0
9	1		0	0
10	0		0	0
11	1		0	0
12	1		0	0
13	1		0	0
14	1		0	0
15	1		0	0
16	1		0	0
17	1		0	0
18	0		0	0
19	0		0	0
20	0		0	0
21	0		0	0
22	0		0	0

31000

Fig. 31

【 図 3 2 】

$freq(C_k, P_j, X): P_j = \text{ノード}3$

X	C _k				
	6	7	8	9	10
6			1	1	1
7	2	2	1	1	1
8	1	1	1	1	0
9	1	1	1	1	0
10	1	1	0	0	
11	1	1		1	0
12	1	1		1	0
13	1	1		1	0
14	1	1		1	0
15	1	1	1		0
16	1	1	1		0
17	1	1	1		0
18	1	1	0	0	
19	1	1	0	0	
20	1	1	0	0	
21	1	1	0	0	
22	1	1	0	0	

32000

Fig. 32

【 図 3 3 】

$freq(C_k, P_j, X): P_j = \text{ノード}8$

X	C _k			
	11	12	13	14
11		1	1	1
12	1		1	1
13	1	1		1
14	1	1	1	

33000

Fig. 33

【 図 3 4 】

$freq(C_k, P_j, X): P_j = \text{ノード}9$

X	C _k		
	15	16	17
15		1	1
16	1		1
17	1	1	

34000

Fig. 34

【 図 3 5 】

$freq(C_k, P_j, X): P_j = \text{ノード}10$

X	C _k				
	18	19	20	21	22
18		1	1	1	1
19	1		1	1	1
20	1	1		1	1
21	1	1	1		1
22	1	1	1	1	

35000

Fig. 35

【 図 3 6 】

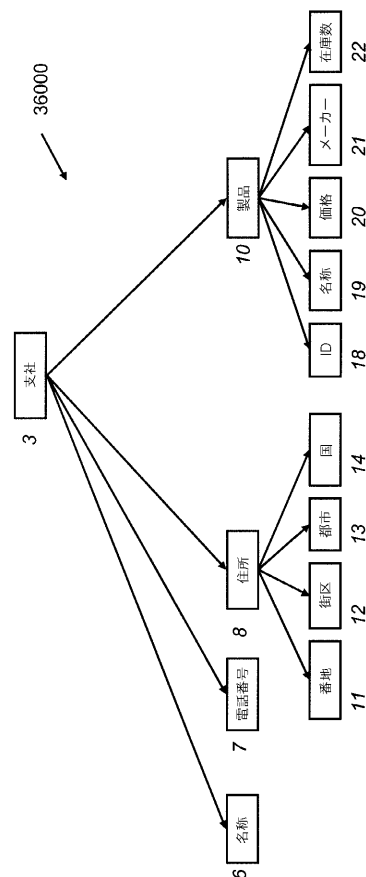


Fig. 36

36000

【図 37】

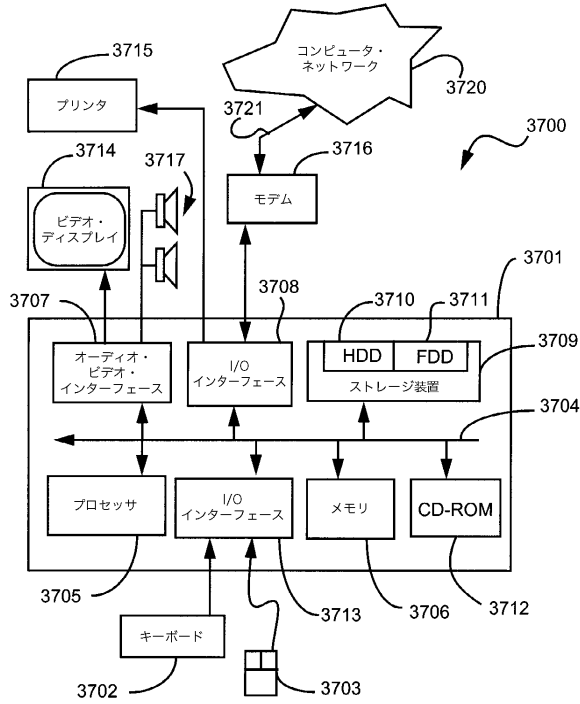


Fig. 37

【図 38】

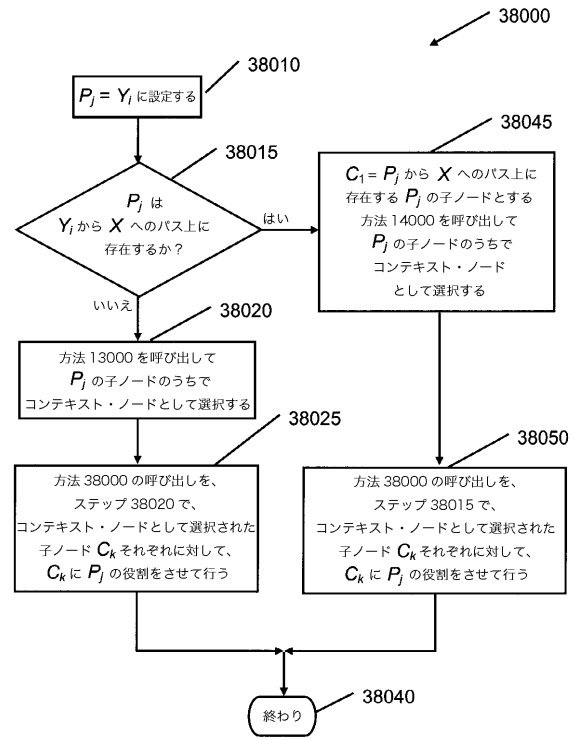


Fig. 38

【図 39】

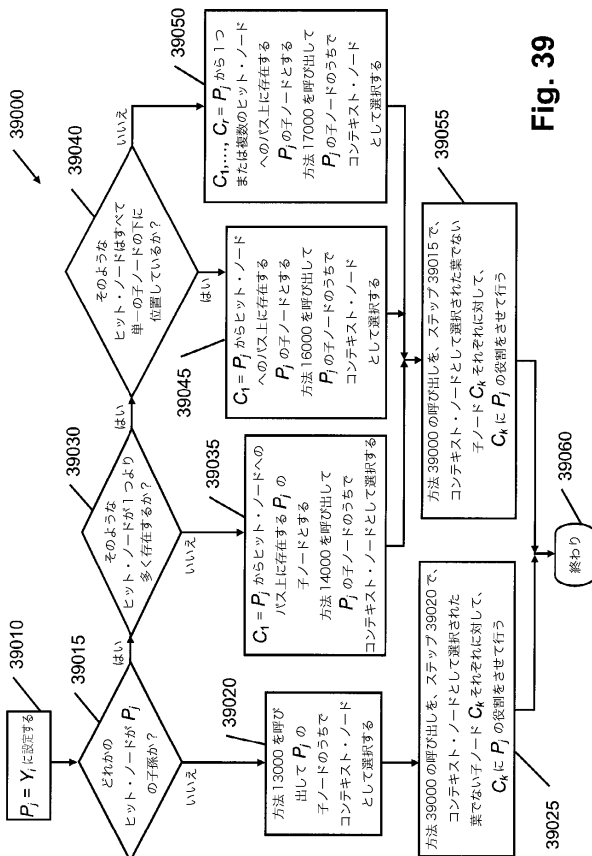


Fig. 39

フロントページの続き

- (72)発明者 ドーン,カーン,フィ,ヴァン
オーストラリア国 2113 ニュー サウス ウェールズ州, ノース ライド, トーマス ホ
ルト ドライブ 1 キヤノン インフォメーション システムズ リサーチ オーストラリア
プロプライエタリー リミテッド 内
- (72)発明者 レノン,アリソン,ジョアン
オーストラリア国 2113 ニュー サウス ウェールズ州, ノース ライド, トーマス ホ
ルト ドライブ 1 キヤノン インフォメーション システムズ リサーチ オーストラリア
プロプライエタリー リミテッド 内

審査官 紀田 馨

- (56)参考文献 特開2003-167879(JP,A)
国際公開第02/027544(WO,A1)
国際公開第01/024045(WO,A1)
池田 ,基礎講座,日経コンピュータ ,日本,日経BP社 Nikkei Business Publications,
Inc.,1998年 3月 2日,no.438,第222頁乃至第227頁

- (58)調査した分野(Int.Cl.,DB名)
JSTPlus/JMEDPlus/JST7580(JDreamII)
G06F 17/30