

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6431055号
(P6431055)

(45) 発行日 平成30年11月28日(2018.11.28)

(24) 登録日 平成30年11月9日(2018.11.9)

(51) Int.Cl. F I
G06F 17/30 (2006.01) G06F 17/30 310A

請求項の数 17 (全 21 頁)

<p>(21) 出願番号 特願2016-521534 (P2016-521534) (86) (22) 出願日 平成26年6月18日 (2014. 6. 18) (65) 公表番号 特表2016-524766 (P2016-524766A) (43) 公表日 平成28年8月18日 (2016. 8. 18) (86) 国際出願番号 PCT/US2014/042888 (87) 国際公開番号 W02014/205046 (87) 国際公開日 平成26年12月24日 (2014. 12. 24) 審査請求日 平成29年6月9日 (2017. 6. 9) (31) 優先権主張番号 61/836, 407 (32) 優先日 平成25年6月18日 (2013. 6. 18) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 315007628 コピーライト クリアランス センター、 インク。 アメリカ合衆国 マサチューセッツ州 O 1923, ダンバーズ, 222 ローズウ ッド ドライブ (74) 代理人 110000659 特許業務法人広江アソシエイツ特許事務所 (72) 発明者 マーマニス, バピス アメリカ合衆国 マサチューセッツ州 O 2451, ウォルサム, 27 シェリル レーン</p>
--	--

最終頁に続く

(54) 【発明の名称】 文献のテキストマイニングのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

ユーザによるアクセスに多様な費用を伴う複数の調査文献の、前記ユーザによるテキストマイニングを容易にするためのシステムであって、

(a) 前記複数の調査文献を格納するように構成されたコンテンツ収納庫であって、前記コンテンツ収納庫は、前記ユーザからの問合せを受けて、テキストマイニングのための前記複数の調査文献の予備的な集合を選択し、前記コンテンツ収納庫は、ユーザが選択的に問合せを変更して、前記ユーザに向けて最適化された前記複数の調査文献の最終的な集合を与えることを可能にする、前記予備的な集合の中の調査文献に関するコンテンツ展開メトリックを提供する、コンテンツ収納庫と、

(b) 導出テキストマイニングデータセットを提供する、調査文献の前記最終的な集合のテキストマイニングのためのテキストマイニングプロセッサと、
 を備えるシステム。

【請求項 2】

前記複数の調査文献のテキストマイニングを管理するプロジェクトマネージャであって、前記プロジェクトマネージャは前記コンテンツ収納庫およびテキストマイニングプロセッサと電氣的にリンクしている、プロジェクトマネージャをさらに備える、請求項 1 に請求されるシステム。

【請求項 3】

前記プロジェクトマネージャが、前記ユーザによる前記システムへの直接のアクセスの

ためのコンピュータインターフェースを提供する、請求項 2 に請求されたシステム。

【請求項 4】

前記コンテンツ収納庫が、集められる前記調査文献の前記コンテンツ展開メトリックに関する 1 つまたは複数のルールに従って問合せを実行する、請求項 3 に請求されるシステム。

【請求項 5】

前記コンテンツ収納庫が、前記予備的な集合の前記調査文献のコンテンツ展開メトリックに関するレポートを生成する、請求項 4 に請求されるシステム。

【請求項 6】

前記レポートが、リスト、円グラフ、線グラフ、および単一の値で構成されるグループからの少なくとも 1 つの表示を含む、請求項 5 に請求されるシステム。

10

【請求項 7】

前記コンテンツ収納庫が、

(a) 前記複数の調査文献のそれぞれのための書誌的メタデータおよびフルテキストを格納するデータ格納デバイスと、

(b) 前記問合せを受け、実行するコンテンツ選択装置であって、前記コンテンツ選択装置が前記データ格納デバイスと電気的に接続されている、コンテンツ選択装置とを備える、請求項 4 に請求されるシステム。

【請求項 8】

前記データ格納デバイスが、前記コンテンツ収納庫が前記ユーザに対し、前記複数の調査文献のそれぞれのアクセス費用を既定することを可能にするユーザアクセス権のデータベースを含む、請求項 7 に請求されるシステム。

20

【請求項 9】

前記コンテンツ選択装置が、文献アクセス費用のパラメータを前記問合せでサポートすることが可能である、請求項 8 に請求されるシステム。

【請求項 10】

前記コンテンツ選択装置が、前記問合せのための前記費用パラメータ内の、前記複数の調査文献のそれぞれのユーザアクセス費用を利用する、請求項 9 に請求されるシステム。

【請求項 11】

前記コンテンツ選択装置が、文献アクセス費用のパラメータを、前記ユーザによって既定され変更可能である前記問合せでサポートすることが可能である、請求項 10 に請求されるシステム。

30

【請求項 12】

前記コンテンツ選択装置が最大ユーザアクセス費用を前記問合せでサポートする、請求項 11 に請求されるシステム。

【請求項 13】

前記コンテンツ収納庫が、前記調査文献の前記予備的な集合の中の調査文献のために、レポートの費用に関するコンテンツ展開メトリックを提供する、請求項 5 に請求されるシステム。

【請求項 14】

前記コンテンツ選択装置が、将来の容易なテキストマイニング操作のための問合せに関連して引き出された、調査文献の最終的な集合を相互参照および格納する、請求項 1 に請求されるシステム。

40

【請求項 15】

前記テキストマイニングプロセッサが、同様のデータ構造の並行クラスタを利用して、前記調査文献の前記最終的な集合のテキストマイニングを実行する、請求項 1 に請求されるシステム。

【請求項 16】

前記テキストマイニングプロセッサが、標準およびカスタムの両方のテキストマイニングプロセッシングモジュールを構築するアプリケーションプログラミングインターフェー

50

スを含む、請求項 15 に請求されるシステム。

【請求項 17】

前記テキストマイニングプロセッサと接続された導出データ収蔵庫であって、前記前記導出データ収蔵庫は前記導出テキストマイニングデータセットを格納する、導出データ収蔵庫をさらに備える、請求項 1 に請求されるシステム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は概して、科学、技術、および医学の分野において発行された調査文献に関し、より詳細には、包括的であるが効率的な方式での調査文献テキストマイニングのためのシステムおよび方法に関する。

10

【背景技術】

【0002】

毎年、数千万もの学術文献が世界中で発行される。それらの発行された文献、または記事の大部分が、無料で提供される特定の記事へのアクセス、および各文献の権利を有する主体により示された費用で提供される他の記事へのアクセスにより、調査者によるレビューに電子的に利用可能である。

【0003】

特定の調査トピックにおいて電子的に利用可能な情報が大量にあるため、調査者にとって、絶えず増え続けるその主題の電子情報を、包括的であるが効率的に検索することがしばしば困難である。特に、調査文献の検索に使用するには、在来の検索エンジンは適していないことがわかっている。この理由は、とりわけ、選択基準の仕様と処理が、関連性について少数の文献を評価するには効果的であるが、すべてが極めて特定の基準にフィットする大量の文献の中から選択する目的には不相当であるためである。このため、特定の主題に関して電子的に利用可能な膨大な量の情報があまりに大量であるため、調査者はしばしば、適切な文献を探し出すことに失敗するリスクがあり、このことは極めて望ましくないものである。

20

【0004】

したがって、調査者による大量の発行された記事の検索を補助するために、組織（例えば、発行者および権利のマネジメントサービス）が、当該技術分野で「テキストマイニング」として知られているプロセスを通して、調査文献のテキストから高品質のデータを解析し、抜粋することを可能にするソフトウェアおよびデータベースを構築することがますます慣習的になってきている。テキストマイニングプロセスの解析を通して、数百万もの文献からのテキストを分析し、相互参照することで、コンピュータベースの検索ツールを使用することにより、調査者が適切な発行物をより効率的に特定すること可能にする。

30

【0005】

発行された調査文献の効率的なテキストマイニングのプロセスには、多くの課題があり、現在、一定の制限がされている。

【0006】

第 1 の課題として、発行された調査文献の効率的なテキストマイニングはまず、関連する文献の全集を大量に集める必要がある。具体的には、より包括的にするために、科学調査のテキストマイニングには、可能なかぎり多くの調査文献にアクセスすることが必要である。同時に、調査文献の集合の権利者は、許可されていない記事の複写および拡散のリスクがあり、このために、権利者が、購読および他の購入されたアクセスの慣習的な形態を通して文献から潜在的に収入を得ることが妨げられることから、テキストマイニングの目的で文献にアクセスすることに同意することをしばしば躊躇する。あらゆる許可されていない文献の複写によるリスクを制限するため、発行者はしばしば、テキストマイニングの目的のために、記事を無作為に（例えば、アルファベット順に並べられた文または単語で）抽出する形式で提供する。しかし、無作為に記事を抽出する形式では、テキストマイニングの特定の機能（例えば、識別された記述パターンに基づいて調査書類と実験記録と

40

50

を区別する能力)が制限されることがわかっている。このため、このやり方は理想的ではないことがわかっている。

【0007】

第2の課題として、発行された調査文献のテキストマイニングには、現在、エンドユーザへの費用を含むことが考慮されていない。上述したように、異なる文献にアクセスするのに別々の費用を伴う。このため、検索経費が制限されている調査者は、無料の発行物の検索に制限することを選択する場合があります。このため、適切な文献を探し出すことに関するリスクがある。同様に、文献にアクセスするのに料金が必要である発行物を含む多くの発行物に検索範囲を広げることを選択した、検索経費が制限されている調査者は、しばしば、過大で法外な調査費用に苦しめられる。

10

【0008】

第3の課題としては、発行された調査文献の効率的なテキストマイニングには、検索結果がエンドユーザに、文献の大きな集団のテキストのすべてへのアクセスを提供することが必要である。対照的に、在来の検索エンジンは、人による評価のための制限された文脈上の情報を伴う、個別の記事へのリンクのリストのみを返す。このことは、調査者が各記事の関連性を判断するには不十分であることがわかっている。

【0009】

第4の課題としては、発行された調査文献のテキストマイニングは現在、エンドユーザに、検索結果に関する有用な問合せ情報を何も提供していない。さらに、エンドユーザは概して、なぜ最初の検索の際に特定の文献が検索されたのかを判断するための制限されたデータを有する。このように、エンドユーザは、以前の検索からの情報を使用して将来の検索の全体の効率性を上げることが妨げられる。

20

【発明の概要】

【発明が解決しようとする課題】

【0010】

本発明の目的は、調査文献のテキストマイニングのための新規の改善されたシステムおよび方法を提供することである。

【0011】

本発明の別の目的は、包括的で費用効率の高い方式での、調査文献のテキストマイニングのためのシステムおよび方法を提供することである。

30

【課題を解決するための手段】

【0012】

したがって、本発明の特徴の1つとして、ユーザによるアクセスに多様な費用を伴う複数の調査文献の、ユーザによるテキストマイニングを簡略化するためのシステムであって、(a)複数の調査文献を格納するように構成されたコンテンツ収納庫であって、コンテンツ収納庫は、ユーザからの問合せを受けて、テキストマイニングに向けて複数の調査文献の予備的な集合を選択し、コンテンツ収納庫は、ユーザが選択的に問合せを変更して、ユーザに向けて最適化された複数の調査文献の最終的な集合を与えることを可能にする、予備的な集合の中の調査文献に関するコンテンツ展開メトリックを提供する、コンテンツ収納庫と、(b)得られたテキストマイニングデータセットを提供する、調査文献の最終的な集合のテキストマイニングのためのテキストマイニングプロセッサと、を備えるシステムが提供される。

40

【0013】

他の様々な特徴および利点が以下の説明から明らかになる。説明中、添付図面が参照され説明の一部を形成する。その中で、本発明を実施するための実施形態が例として示される。実施形態は、当業者が本発明を実施することを可能にするように十分に詳細に記載され、他の実施形態が利用でき、本発明の範囲を逸脱することなくその構造の変更がされることが理解される。したがって、以下の詳細な説明は限定の意味としてとられず、本発明の範囲は、添付の特許請求の範囲により最適に既定される。

【0014】

50

図中、同様の参照符号が同様の部品を示す。

【図面の簡単な説明】

【0015】

【図1】図1は、文献のテキストマイニングのためのシステムの簡略化されたブロック図であり、このシステムは本発明の教示にしたがって構成されている。

【図2】図2は、図1に示すコンテンツ収納庫に格納された、記事に関する様々な形式のデータ間の実施可能な関係を理解するのに有用な、例示的なデータモデルである。

【図3】図3は、図1に示す文献収納庫内の記事のアクセスドメインの実施を理解するのに有用な、例示的なデータモデルである。

【図4】図4は、図1に示すシステムを使用する文献のテキストマイニングの新規の方法の簡略化されたフローチャートである。

【図5】図5は、図4に示すテキストマイニングの方法のより詳細なフローチャートである。

【図6】図6は、図1に示すコンテンツ選択装置に格納された、展開メトリックに関するデータの実施可能な関係を理解するのに有用な例示的なデータモデルを示す。

【図7a】図1に示すシステムの例示的な使用を理解するのに有用な、一連のサンプルスクリーン表示である。

【図7b】図1に示すシステムの例示的な使用を理解するのに有用な、一連のサンプルスクリーン表示である。

【図7c】図1に示すシステムの例示的な使用を理解するのに有用な、一連のサンプルスクリーン表示である。

【図7d】図1に示すシステムの例示的な使用を理解するのに有用な、一連のサンプルスクリーン表示である。

【図7e】図1に示すシステムの例示的な使用を理解するのに有用な、一連のサンプルスクリーン表示である。

【発明を実施するための形態】

【0016】

テキストマイニングシステム11

ここで、図1を参照すると、調査文献のテキストマイニングのためのシステムの全体的なブロック図が示されている。このシステムは本発明の教示にしたがって構成され、全体として参照符号11で識別されている。以下でさらに詳細に説明されるように、システム11は、とりわけ、(i)それに続くテキストマイニング操作の主題となる、費用パラメータを調査文献の集合の選択のプロセスに組み込み、それによって、(ii)選択された文献の展開に関する、ユーザが直観的に認識可能なメトリックを提供するように設計される。必要であれば、ユーザは、この後に、テキストマイニングされる調査文献の最適化された集合をもたらすために、メトリックを利用して文献選択プロセスの特定のパラメータを修正することができる。この能力で、システム11は、包括的であるが費用効率がある、調査文献の集合のテキストマイニングを促進する。このことは、本発明の主要な目的である。

【0017】

例示の目的のみのために、システム11は本明細書において、調査文献の大きな収納庫を使用して行われるテキストマイニング作業と関連づけて説明される。しかし、システム11は調査文献のテキストマイニングに限定されないことが理解される。さらに、システム11は、あらゆる形式の文献、特に、アクセスするのに費用を伴うあらゆる文献から関連するテキストを識別することを要するあらゆる環境で使用可能であることが理解される。

【0018】

システム11は複数のモジュールを含み、この複数のモジュールは全体で、本発明のテキストマイニング操作をエンドユーザ13に提供する。具体的には、以下に詳細に記載されるように、システム11は、システム11の中心の機能的ハブとして働くプロジェクト

10

20

30

40

50

マネージャ 15 と、テキストマイニングおよび計量されたアクセスのための記事を含む文献収納庫 17 と、本発明の主要なテキストマイニング作業を行うテキストマイニングプロセッサ 19 と、テキストマイニングプロセッサ 19 によって行われるテキストマイニング操作の出力を格納する導出データ収納庫 21 とを備える。

【0019】

プロジェクトマネージャ 15 は本明細書において、任意の通信媒体を介して（例えば、インターネットを介して）、エンドユーザ 13 のコンピュータデバイスと電子的にリンクしているサーバとして示される。この方式で、プロジェクトマネージャ 15 は、エンドユーザ 13 に、システム 11 にアクセスするための主要なインターフェイスを提供する。以下にさらに詳細に記載されるように、プロジェクトマネージャ 15 はエンドユーザ 13 に、（i）新たなテキストマイニングプロジェクトを構築すること、（ii）進行中のプロジェクトの状態および進捗を追跡すること、および、（iii）完了したプロジェクトから返されるデータにアクセスすることを可能にする。

10

【0020】

なお、テキストマイニングプロジェクトへのアクセスは、個別、団体ベース、および組織のいずれかのレベルのアクセス権で、プロジェクトマネージャ 15 から所与のエンドユーザ 13 に与えることができる。この能力で、システム 11 が様々な異なる環境で実施され得ることが見込まれる。

【0021】

文献、またはコンテンツ収納庫 17 は、学術記事の大きい集まりの書誌的メタデータとフルテキストとの両方を含むデータ格納デバイス 23 - 1 および 23 - 2 を備え、コンテンツは、速やかに引き出すことを容易にするように、索引がつけられていることが好ましい。

20

【0022】

たとえば、ここで図 2 を参照すると、コンテンツ収納庫 17 に格納された、記事に関する様々な形式のデータ間の実施可能な関係を理解するのに有用な、例示的なデータモデルが示されている。このデータモデルは全体として、参照符号 25 で識別される。しかし、本発明の趣旨を逸脱することなく、他のデータベース技術の類似のデータモデルが同様に、データベースモデリングの経験を積んだ専門家によって構築され得ることが理解される。

30

【0023】

見て取れるように、データモデル 25 は、各記事のメタデータを有する記事テーブル 27 を含み、このデータは、それに限定されないが、ワークのタイトル、ワークの著作者、および特定のキーワードを含む。記事テーブル 27 は、各記事のフルテキスト（すなわち、発行された形式の文献を構成する完全な原本に基づく事項）、ならびに、書誌的事項、引例のリスト、および/または、収納庫 19 に置かれているか、もしくは置かれていない別の記事のセットの参照をさらに含むことが好ましい。

【0024】

著作者テーブル 29 は、（記事著作者テーブル 31 を介して）記事テーブル 27 とリンクしており、学術文献を創作する様々な個人または組織を表す。著作者は文献収納庫 17 に、名前によって、また、基本的な識別名の任意のセットをともなって出されることが好ましい。

40

【0025】

出所テーブル 33 が、記事の一般的なソースに関するデータ（すなわち、どこで記事を見つけられるか）を提供する。ジャーナル（すなわち、記事のセットを発行する学術ワーク）と収蔵庫とは両方とも出所の形式である。したがって、ジャーナルテーブル 35 が、その中に見られるタイトル、標準番号、および発行者を含む各ジャーナルの属性とともに、出所テーブル 33 とリンクする。同様に、コレクションテーブル 37 が出所テーブル 33 とリンクし、ジャーナルと集合との両方に潜在的に見られる記事とともに、記事の代替的なソースを提供する。

50

【 0 0 2 6 】

最後に、発行テーブル 3 9 が記事テーブル 2 7 および出所テーブル 3 3 内のデータ間の関係を構築する。発行テーブル 3 9 は、発行者から直接、しばしば高額での記事の利用可能性を示すデータを含む。たとえば、特定の記事は、そのももとの発行者から 4 0 ドルで、また、ドキュメントの収納庫から 5 ドルで、利用可能である場合がある。

【 0 0 2 7 】

したがって、例示的なデータモデル 2 5 の構造を使用することで、他の事を含め、(i) 著作者もしくは著作者のセット、(i i) 記事のタイトル、(i i i) キーワードもしくは他の同様のメタデータの領域、(i v) 発行物もしくは発行物のセット、(v) ジャーナルもしくはジャーナルのセット、(i v) 集合もしくは集合のセット、および/または (v i i) 発行データの範囲に関するデータを使用して検索の問合せを容易に処理できることが明らかである。

【 0 0 2 8 】

少なくとも 1 つのデータ格納デバイス 2 3 が、ユーザのアクセス権のデータベースをさらに含むことが理解される。したがって、権利の付与、またそれによって、問合せ、ジョブ、およびユーザによる記事レベルのアクセス記録に基づき、文献収納庫 1 7 は各ユーザのアクセス権を追跡することができる。

【 0 0 2 9 】

たとえば、図 3 を参照すると、文献収納庫 1 7 内の記事のアクセストメインの実施を提供する例示的なデータモデルを示しているこのデータモデルは全体として参照符号 4 1 で識別される。見て取れるように、データモデル 4 1 は、エンドユーザテーブル 4 3 と団体テーブル 4 5 とを (団体ユーザテーブル 4 7 を介して) 相互参照する。この理由は、各団体は通常、複数の異なるユーザを含むためである。さらに、団体はしばしば、購読を複数購入するため、団体テーブル 4 5 は購読テーブル 4 9 とリンクしている。記事のソース (たとえば、その中の記事が購入可能な異なる集合) を明示する出所テーブル 5 1 は、そのため、購読アイテムテーブル 5 3 を介して購読テーブル 4 9 とリンクしている。したがって、システム 1 1 は、エンドユーザ 1 3 に文献収納庫 1 7 に含まれる大量の記事を効率的にテキストマイニングすることを可能にするのみならず、各エンドユーザ 1 3 がどの記事を購入しているかを容易に確認する。このことは、非常に望ましいことである。

【 0 0 3 0 】

図 1 に戻ると、文献収納庫 1 7 は、データ格納デバイス 2 3 とプロジェクトマネージャ 1 5 との両方に接続された、コンテンツ選択装置、または問合せプロセッサ 5 5 をさらに有する。したがって、以下にさらに記載されるように、コンテンツ選択装置 5 5 はデータ格納デバイス 2 3 から調査文献にアクセスし、様々な異なるフルテキストおよびメタデータの問合せを実施することにより、記事の、最適化されたサブセット、またはクラスタを選択する。結果として得られた文献のクラスタは、その後、将来の問合せを容易にするために、コンテンツ選択装置 5 5 に格納され、文献のクラスタは、元の問合せが繰り返された際に、必要に応じてアップデートされる。

【 0 0 3 1 】

本発明の主要な特徴としては、コンテンツ選択装置 5 5 が、データ格納デバイス 2 3 から最初の文献の集団を与えるために、費用パラメータをフルテキストおよびメタデータの問合せに含むことできる。さらに、コンテンツ選択装置 5 5 は、エンドユーザ 1 3 に、最初の問合せから得られた選択された文献の展開に関する直観的メトリックを与える。この方式で、ユーザは、以下にさらに説明するように、続いてテキストマイニングされる調査文献の、包括的であるが効率的な展開を得るために、必要に応じて問合せを改善できる。

【 0 0 3 2 】

上記で簡潔に参照したように、テキストマイニングプロセッサ 1 9 は、本発明の主要なテキストマイニング作業を担っている。言い換えると、テキストマイニングプロセッサ 1 9 は調査者に、収納庫 1 9 から引き出された文献の関連する集合のテキストマイニングジョブを特定することを可能にし、そのジョブを、そのジョブの要求と非同期的に実行し、

10

20

30

40

50

そして、完了した後に調査者に知らせる。

【0033】

本明細書で述べるように、テキストマイニングプロセッサ19は、標準化されたアーキテクチャにしたがって、テキストマイニングプログラムを並行して実行するように設計された、スタックされた複数の計算装置57-1~57-3を備える。具体的には、テキストマイニングソフトウェアが入力データを導出データ収納庫21内の計算装置59-1~59-3から受け取り(すなわち、前回のテキストマイニング操作の出力)、文献のセット内に特定された集合に向けて、文献のメタデータおよびフルテキストのテキストマイニング操作を並行して実施して、その後、導出データ収納庫21内の名前が付されたデータセットに格納される出力を与える。各作業に向けた処理資源の分配は、テキストマイニングプロセッサ19によって内部で追跡されることが好ましい。

10

【0034】

システム11を使用するテキストマイニング法111

上記で簡潔に参照したように、システム11は調査文献のテキストマイニングの新規の方法に関わるように設計される。具体的には、ここで図4および5を参照すると、システム11を使用するテキストマイニングのための文献の選択、購入、および処理の新規の方法の、簡略化されたフローチャート、およびわずかにより詳細なフローチャートがそれぞれ示されている。この方法は本明細書において、全体として参照符号111で識別されている。

20

【0035】

以下にさらに詳細に説明するように、本発明のテキストマイニング法は、最初に、検索変数、またはパラメータ値のセットを使用して調査文献の集団、またはプールを集めて、潜在的に関連する調査文献の広範囲の集合を与える。言い換えると、最初の集合は、問合せの基準に最も合う単一の文献を見つけることを試みるような、人間の選択のために、関連性により優先順位をつけられた文献を返すことを求めている。代わりに、結果のセットは審査のために提供されるのではなく、むしろそれに続くテキストマイニング処理のために集められる。

【0036】

前述の文献選択プロセスは、複数の記事の回りに「フェンス」を投じて、集合のサブセットを形成することに類似している。そのため、このフェンスの構成は、後に、コンテンツの広げられたメトリック(たとえば、なぜ特定の文献が最初に選択されたかについての情報)を使用して、ユーザによって変更され、エンドユーザ13にとって最も適切で望ましい選択(たとえば、費用、発行者など)に、元の調査文献のプールを再び既定、または狭めることができる。この方式で、それらのすべてが特定の性質に従う調査文献の高品質の選択が、それに続く、効率的かつ費用効率のよい方式のテキストマイニング操作のために集められる。

30

【0037】

なお、テキストマイニングジョブは、プロジェクトマネージャ15にアップロードされるプログラムコードで構成される。

【0038】

プロセス111を開始するために、エンドユーザ13はまず、テキストマイニングプロジェクトを既定、または構築する。このプロジェクト既定ステップは、全体として参照符号113で識別される。具体的には、プロジェクト既定ステップ113の一部として、エンドユーザ13は、(i)テキストマイニング操作に利用される文献セット(すなわち、収蔵庫19内のコンテンツの選択)、(ii)プロセスの明細(すなわち、文献のトークン分解、特有の属性の計算、および同様のデータ構造の並行クラスタリング)、ならびに(iii)明細の報告(すなわち、ユーザにテキストマイニングの結果を提供するための特定の手段)を特定する。

40

【0039】

なお、文献セットは、(i)文献の識別名、著作者、共同制作者、機関、および発行者

50

など（もしくは、上述の属性の任意のリストまたは集合）の明細を使用する文献の問合せを通して、または（i i）あらかじめ既定された文献セット（すなわち、以前の問合せの結果として得られた文献セット）を使用することで、特定することができる。

【0040】

ステップ113が完了すると、コンテンツ選択装置55は、ステップ113で特定された、あらゆるコンテンツ展開条件（たとえば、「C. Elegans」という用語を含むが、発行者Xによる記事は除くすべての文献を特定する）を受けて、ジョブに向けて調査文献を選択する。この文献選択ステップは、全体として参照符号115で識別されている。

【0041】

文献選択ステップ115の一部として、システム11は、エンドユーザ13が、最初の文献の集合に関連する展開メトリックを識別し、分析することを可能にするユーザインターフェイスを生成する。この能力で、エンドユーザ13は予備的な問合せの特定のパラメータを変更して、テキストマイニングされる、より最適な文献の集合を与えることができる。

【0042】

対照的に、在来のテキストベースの検索結果は通常、説明がされない。言い換えると、ユーザは概して、なぜ検索結果が特定され、特定の順番で並べられるのかを理解しない。しかし、調査フィールドでは、調査者は検索要求からの任意のコンテンツの選択を利用できない。大量の調査記事が利用可能であるため、調査者は、なぜ特定の記事が選択されたのか、より重要には、どのように検索パラメータの重要性、または詳細を変更して検索結果に影響を与えたのかを知る必要がある。

【0043】

したがって、上記で簡潔に参照したように、問合せプロセッサ55は、選択された検索メトリック（すなわち、コンテンツ、発行者、費用などによる検索結果の内訳）に基づき、ユーザに向けてレポートを生成する。この方式で、エンドユーザ13は検索結果に影響を与えた要因をよりよく特定することができる。したがって、システム11はエンドユーザ13が、この後、作動中に（on the fly）検索パラメータを調整し、それに続く文献の第二の集合を実施して、前回の集合内に発見されたあらゆる不十分さを順応させることを可能にする。

【0044】

ステップ115で最初に集められた調査文献の集団が拡大したことにより、文献処理ステップが、その中の文献の最適化されたグループ、またはサブセット（すなわち、識別された特定のキーワードに関して最も近似の文献）を既定、または識別し始める。この文献処理ステップは全体として参照符号117で識別される。

【0045】

文献処理ステップ117は好ましくは、大きなデータセットのバッチ処理に使用される、パイプライン方式のマップの様々な縮小パラダイムを利用する。好ましくは、テキストマイニングプロセッサ19は、カスタムマップの構築のためのアプリケーションプログラミングインターフェイス（API）を提供し、モジュールを削減する。

【0046】

具体的には、個別の文献の操作を行って各文献を他の形態に変換する「マップ」処理が特定され得る。たとえば、プロセスにより、遺伝子配列リサーチを記載する文書が、各文書によって言及される特定の遺伝子のリストに変換され得る。

【0047】

さらに、「削減」プロセスは、変換された文書のリストを集合の形態に結合する。たとえば、プロセスにより、調査文書の集合によって言及される遺伝子のリストを取り、したがって、調査を行う機関により集計された遺伝子のリストを返すことができる。削減変換の第二の段階が、第一の段階の出力を操作して、機関による遺伝子のセットを取り、機関による集計を繰り返す。これは「結合」変換と呼ばれる。このやり方でプロセスを分割す

10

20

30

40

50

ることで、ジョブを並行して実施することを補助するのに役立つ。

【0048】

本発明の新規の特徴として、文献処理ステップ117は、標準処理モジュール119とカスタム処理モジュール121との両方を補助し、以下にさらに説明されるように、それらからの出力はさらに処理されて独自の属性を見出す。

【0049】

標準処理モジュール119は、すべてのエンドユーザ13が使用するために、テキストマイニングプロセッサ19によって提供される。標準処理モジュール119の例は、調査作業の専門性を高めるために、(i)記事の、セクション、段落、文、および単語の階層へのトークン分解(すなわち、分解または分裂)、(ii)スピーチのタグ付け(すなわち、単語を名詞や動詞などとして識別すること)の一部、(iii)引例の抜粋(すなわち、記事の書誌的事項を記事のメタデータまたは記事の参照のリストに変換すること)、ならびに(iv)因子の抜粋(すなわち、HOXA1、BRCA1などのHUGO gene nomenclature systemによる、記事の単語形式のタグ付け)を含む。

10

【0050】

カスタム処理モジュール121は、繰り返し使用するために、特定のエンドユーザ13により構築され、モジュールのアプリケーションプログラミングインターフェイス(API)に係るプログラムとして実施される。本発明の特徴として、カスタム処理モジュール121は、その構築を担うエンドユーザにより個人的使用のために保持されるか、匿名または記名形式ですべてのエンドユーザ13による広範囲の使用のために発行される。多くの顧客に頻繁に利用されるカスタム処理モジュール121により、その創作者に特別な特典または金銭的利益を与えられる場合があることが理解される。

20

【0051】

文献の最初の集合がテキストマイニング処理モジュール119および121により、分解され、タグが付され、および/または変換されると、次いで、独特の、ユーザにより特定された属性が識別されてデータセット123を形成する。次いで、データセット123は、以下にさらに詳細に説明されるように、データ削減、または関連データを並行してクラスタにする集合処理ステップ125の際に更に削減される。

【0052】

データ削減ステップ125は、標準データセット処理モジュール127およびカスタムデータセット処理モジュール129へのアクセスにより、それぞれモジュール119および121を増大して、標準データセットおよびカスタムデータセットを与える。

30

【0053】

標準データセットは、対になっているデータ(すなわち、名前および数値)の集合であり、したがって、任意のモジュールから名前でアクセスすることができる。標準データセットの例には、それに限定されないが、ISO国名コード、HUGO gene nomenclature、および元素周期表が含まれる。

【0054】

カスタムデータセットは標準データセットに似ているが、システム11の個別のエンドユーザ13により与えられる。カスタムモジュールのように、カスタムデータセットは個人的使用のために保持されるか、または、匿名か記名のいずれかでシステムのすべてのエンドユーザ13による使用のために発行される。再度、多くの顧客に頻繁に利用されるカスタムデータセットは、その創作者に特別な特典または金銭的利益を与える場合があることが理解される。

40

【0055】

データセット処理モジュール127および129は、パイプライン、またはクラスタに結合される。モジュール127および129の出力は、別のデータセット処理モジュールに直接流れることができる。または、いくつかのデータセット処理モジュールの出力が集合およびフィルタリング操作を使用して結合されることができる。

50

【 0 0 5 6 】

ステップ 1 2 5 における関連するデータの並列クラスタリングが完了すると、報告ステップ 1 3 1 の一部として、テキストマイニング操作の結果がユーザ 1 3 に報告される。報告ステップ 1 3 1 において、標準報告モジュール 1 3 3 およびカスタム報告モジュール 1 3 5 が、テキストマイニング操作から最も適切と思われる文献の書誌データを生成する。この書誌データは収納庫 2 1 内に導出データセットとして格納される。この導出データセットはこの後、プロジェクトマネージャ 1 5 を介する調査の過程の中で、エンドユーザ 1 3 によって引き出され、検査することが可能である。

【 0 0 5 7 】

コンテンツ選択装置 5 5 の費用計算モジュール

上記で簡潔に参照したように、コンテンツ選択装置 5 5 は、エンドユーザ 1 3 に、テキストマイニングのために文献の最適な集合を引き出すことを確実にする、相互作用するコンテンツ選択プロセスに従事することを可能にする。本発明の特徴として、コンテンツ選択装置 5 5 は新規の費用計算モジュールを使用したフルテキストおよびメタデータの問合せから引き出された文献の最初の集団を改善または最適化することができる。言い換えると、記事にアクセスする費用を考慮しつつ（すなわち、ユーザがどの記事を購読しているか、検索経費の最大値がいくらか、など）、コンテンツ選択装置 5 5 は、エンドユーザ 1 3 が記事のプールを（たとえば、特定のキーワードに基づいて、記事の言語によって、および/または特定の著作者によって）選択することを可能にするようにプログラムされている。

【 0 0 5 8 】

認識できるように、費用ベースの文献の集合の選択は、調査者にかなりの経済的負担を課し得る。より詳細には、文献収納庫 1 7 は、ユーザ 1 3 が購読していないが、必要なアクセス費用を支払うことで利用可能である非常に多くの記事のテキストを含むか、アクセスできることが好ましい。しかし、在来のテキストマイニングプロセスが通常、エンドユーザに、著作者が読むことを望むか、望み得るよりも多くの文献のアクセスを提供することを考慮すると、的確さが不十分である文献選択の問合せは、実施するには費用が法外になる場合がある。

【 0 0 5 9 】

したがって、コンテンツ選択装置 5 5 には、とりわけ、さらなる検索の制約が存在する一方で、各テキストマイニングジョブのためのコンテンツ費用の最大値を設定し、それを履行するのに使用できる費用計算モジュールが提供される。

【 0 0 6 0 】

テキストマイニングジョブのためのコンテンツ費用の最大値を設定するために、コンテンツ選択装置 2 5 により以下の式が利用され得る。

$$\sum_{d \in D} F(d) \quad (1)$$

式中、 n は集合の中の文献の数、 $F(d)$ は、発行テーブル 3 9 の例示的スキームで定められたように（すなわち、既存の記事の購読/購入を除いて）、各文献 d を得るための費用を定める関数である。

【 0 0 6 1 】

しかし、式 (1) は、ユーザがすでにアクセスする権利が与えられている文献を考慮していない。文献の異なる出所（すなわち、ソース）が異なる平均金額を付けるが、同時に、すべての出所がすべての文献を差し出すわけではないことを考慮することも、有用である。たとえば、文献は (i) ユーザが現在購読をしている出所から費用無しで、(ii) J S T O R (商 標) 電 子 図 書 館 などの公共の文献収蔵庫から低い定額で、および、(iii) 個人の発行者から比較的高額で、使用可能となり得る。したがって、コンテンツ選択装置 5 5 によって利用される、より有用な費用計算式には、以下に示すように、すべての使用可能な出所から取る場合に、各記事のすべての異なる費用の合計を考慮する。

$$\sum_j \sum_{d \in D_j} F_j(d) \quad (2)$$

式中、 n は集合の中の文献の数、 $F(d)$ は、発行テーブル 39 の例示的スキームで定められたように、各出所 j から各文献 d を得るための費用を定める関数である。

【0062】

式(2)を利用して、以下に示すように、問合せセットの制約を足し合わせることで、テキストマイニングジョブのためのコンテンツ費用、または経費(予算)の最大値 B が定められる。

$$\sum_j \sum_{d \in D_j} F_j(d) < B \quad (3)$$

【0063】

任意選択的に、テキストマイニング調査は、例外を減らすし、別の方法で結果の統計上の信頼性を上げるために、選択された調査文献のプールを最大化することを試みる。経費の制約を満たす1つの方法は、同時に、文献の集団を最大化しつつ、費用を上げることで集合の中の記事をソートすることである。記事はその後、順に、記事の集められたセットが既定の経費に達するまで選択される。

【0064】

しかし、上述のような、費用が高くなる選択プロセスを利用することは、特に、文献毎の費用が明らかに異なる多くのプールで構成される多数の文献の場合、多くの調査ジョブの要求に対して極めて不十分である。最も明白には、経費が制約された選択では、無料のコンテンツ、ユーザによって購読されるコンテンツ、および公共の文献収蔵庫のより古いコンテンツが強く重視され、したがって、信頼性および関連性が低い大量の文献を含む検索結果が与えられる。

【0065】

したがって、本発明は、コンテンツの支出の制約を尊重するが、同時に、特定の無料または低価格の出所または他のメタデータのフィールド値への不公平な割当を回避する文献の集団を特定し、選択するための機構を含む。

【0066】

本明細書において既定されるように、「コンテンツ展開」は、文献の集団が、出所によってなど、特定の条件を満たすものの中に広く分配された範囲を示す。たとえば、無料と支払いがされたものとの両方を含み、および様々な異なる発行者からの集合を伴う、多くの異なるソースからの公正な代表を含む調査文献の集団は、比較的広い、または広域のコンテンツ展開を考慮する。

【0067】

コンテンツ選択装置 55 により最初の文献の集合が完了したが、対応するテキストマイニングジョブの実際のスケジューリングおよび実施の前に、コンテンツ選択装置 55 は、様々な予め既定されたメトリック、またはルールを使用してコンテンツ展開を算定する。したがって、コンテンツ選択装置 55 は、1つまたは複数のユーザインターフェイス(UI)のレビュースクリーンを通して、算定されたコンテンツ展開を表示する。この方式で、エンドユーザ 13 は様々なメトリック(たとえば、費用、ソースなど)にわたるコンテンツ展開を分析することができ、必要であれば、テキストマイニング操作のスケジューリングの前に、調整された文献の集合のセットを与えるように、検索パラメータを変更する。

【0068】

コンテンツ展開メトリックは、構成可能な警告の閾値、およびユーザへのメッセージ表示をサポートして、それに続くテキストマイニング操作の際に最適化された文献の集合が利用されることを確実にすることができる。さらに、ユーザは属性、および合計または平均などの集計機能を選択することで、集合の中の文献の様々な異なる属性の中のコンテン

10

20

30

40

50

ツ展開を調査することができる。したがって、コンテンツ選択装置 55 は、セットの要素にわたる集計を計算する。

【0069】

ここで図 6 を参照すると、コンテンツ展開メトリックのそれぞれ、および展開メトリックのルールのそれぞれの結果を実施および表示する手段に関連する、既定、またはルールの柔軟な性質をサポートする例示的なデータモデルが示されている。このデータモデルは、全体として参照符号 211 で識別される。見て取れるように、展開のメトリックテーブル 213 のそれぞれは、複数の変更可能なルール 215 によって既定され、それによりユーザが（閾値テーブル 217 を介して）閾値を使用して展開メトリックを創作して、特定のコンテンツ選択方法に対処することを可能にする。したがって、変更可能なルール 215 のそれぞれは、ユーザに、実施された展開メトリックルールのそれぞれを（たとえば、リスト、円グラフ、線グラフ、および/または単一の数値で）表示する好ましい手段を確立することを可能にする。

10

【0070】

コンテンツ選択装置 55 による展開メトリックルールの利用は、複数ステップのプロセスを要する。プロセスの第 1 のステップでは、適切と考えられれば、変更利用可能であるメトリックに実施される各ルールの規定により、エンドユーザ 13 が、コンテンツ選択プロセスの際に利用される、関連する展開メトリックを選択する。展開メトリックテーブル 213 は好ましくは、エンドユーザ 13 に利用可能なすべての展開メトリックを列挙する。

20

【0071】

特定の展開メトリックの選択がされると、必要であれば、その展開メトリックのための対応する展開メトリックルールが試験および変更のために利用可能になる。展開メトリックルールを既定するための例示的な擬似コードを以下に示す。

```
return true
If count(article) > 1000
    return true
If metric-columns includes - any
(article.author, article.author.institution)
    return true
```

30

【0072】

各展開メトリックテーブル 213 の関連性表示欄は、テキストマイニングジョブの既定に対して実行されて、所与の展開メトリックの関連性について「true」または「false」の値を返すプログラムコードを含む。言い換えると、上記で与えられたルールの第 1 段階に基づいて、「true」の値が、そのルールに関連性がある、またはそのルールが適用されるべきであることを示す。

【0073】

ルールの第 2 段階では、そのルールのパラメータが既定される。提示の例では、1000 より多くの記事がコンテンツ展開の中に存在するかが定められる。このルールは、ジョブの既定に対して実行される集計機能に基づき、関連すると見なされる。

40

【0074】

ルールの第 3 段階では、測定の属性が定められる。その後、前述のプロセスがすべての実行される展開メトリックのルール（すなわち、「true」として識別される関連性の表示を有する各ルールに対して繰り返される。

【0075】

プロセスの第 2 のステップでは、すべての関連する展開メトリック（すなわち、コンテンツ選択プロセスに適用されるメトリック）がコンテンツ選択装置 55 により引き出されることで、それに従って実行される。なお、所与の展開メトリックは、1 つまたは複数の展開メトリックルールを含むことができる。

50

【0076】

このルールの表示欄は、ジョブの既定およびその関連する文献の集合に対して実行されるプログラムコードを含む。例示的な擬似コードを以下に示す。

```
Select article.publication.origin,
       count(distinct article.publication.origin)
       / count(article)
from job.articles
```

【0077】

上記の例示的コードでは、文献のソースのリストが、集団全体におけるその割合によりソートされ、それに応じて表示される。これにより、調査者が、特定のジョブのための文献の集合の中で特定の記事のソースが大きな比率を占めているかを判断することができる。

10

【0078】

さらなる例示的な擬似コードを以下に示す。

```
Select sum(article.publication.price)
from job.articles
```

【0079】

上記の例示的コードでは、特定のジョブに含まれる記事のためのコンテンツ取得費用全体がユーザに向けて表示される。

20

【0080】

このプロセスの最終ステップでは、実施された展開メトリックのそれぞれのリンクが表示されることで、ユーザが展開メトリックルールに示す表示方法にしたがって結果を精査することができる。一例として、円グラフの表示方法により、ルールが、割合として解釈され得る {記事名、記事の値} の対のリストを返すことが示される。別の例として、単一の値の表示方法により、ルールが、メッセージの属性（たとえば、C言語のストリングである「The total cost of the job is %d」、%dのパラメータは表示のために、ルールの表示によって返される値と入れ替えられる）と合わせられ得る単一の値を返すことが示される。

【0081】

上述のジョブの集合のコンテンツ選択プロセスは、強制されたプログラミングまたは最適化技術を使用して達成され得ることが理解される。したがって、知識を有する専門家が、シンプレックス、ミニマックス、および非線形反復法を含む様々な数学上の最適化の方法を利用して、最適に文献収蔵庫19からコンテンツを選択できる。

30

【0082】

テキストマイニングシステム11および方法111の例示的使用

ここで図7(a)-(e)を参照すると、本発明の原理を理解するのに有用である一連のサンプルスクリーン表示が示されている。

【0083】

上記で参照したように、方法111の第1ステップ113は、エンドユーザ13がテキストマイニングジョブを既定することを要する。ステップ115で集められる記事の選択を補助するために、システム11は、コンテンツを選択するためのユーザインターフェイスを生成する。ユーザインターフェイスの例示的なスクリーン表示が図7(a)に示されており、全体として参照符号311で識別される。

40

【0084】

見て取れるように、コンテンツ選択ユーザインターフェイス311は、新しい、または以前に規定されたテキストマイニングプロジェクトへのアクセスを提供する複数のタブ313-1および313-2を含む。各プロジェクトスクリーンは、ジョブを識別するためのプロジェクト名ウィンドウ315、ジョブの範囲を簡潔にまとめるための解説ウィンドウ317、コンテンツ選択プロセスで使用されるキーワードを入力のためのキーワードウ

50

インドウ 3 1 9、コンテンツ選択プロセスから選択された著作者を含むか引き出すための著作者ウインドウ 3 2 1、コンテンツ選択プロセスから選択された発行者を含むか引き出すための発行者ウインドウ 3 2 3、および既定の期間内に発行された記事にコンテンツ選択プロセスを制限するための期日ウインドウ 3 2 5 を含む。ともに、スクリーン 3 1 1 に提供される様々な検索パラメータ、または要素がコンテンツ選択装置 5 5 に渡されてテキストマイニングジョブのために記事の集合を取り込む。

【 0 0 8 5 】

なお、コンテンツ選択ユーザインターフェイス 3 1 1 にはさらに、ユーザが特定のテキストマイニング処理属性を選択および変更することを可能にする属性セットドロップダウンウインドウ 3 2 7 が与えられる。たとえば、ウインドウ 3 2 7 の「value」の語をクリックすると、エンドユーザ 1 3 は、テキストマイニング操作のための検索費用キャップが発効され得る別のスクリーンに移される。

10

【 0 0 8 6 】

具体的には、ここで図 7 (b) を参照すると、コンテンツ展開の範囲を設定するためのユーザインターフェイスのサンプルスクリーン表示が示されている。この例示的スクリーン表示は全体として参照符号 3 3 1 で識別される。見て取れるように、様々な費用に関連するルールが文献選択ステップ 1 1 5 に包含され得る。ユーザインターフェイス 3 3 1 を通して、エンドユーザ 1 3 はリストからルールを選択することで費用の範囲を設定することができ、したがって、ルールのために返された値に対して実行される式を特定する。

【 0 0 8 7 】

20

たとえば、第一ルール 3 3 3 では、式は結果となる最大値が 5 0 になることを示す。言い換えると、記事の集団全体の 5 0 % より多くを構成するソースは無い。コンテンツ選択ステップ 1 1 5 の実行中、コンテンツ選択装置 5 5 は集合のための記事の選択を制限して特定された範囲を受け入れる（すなわち、単一の記事のコンテンツのホットスポットを防止する）。この制限は、したがって、集合に現れる文献数の合計に影響する。

【 0 0 8 8 】

第二のルール 3 3 5 では、式が、このルールによって計算される記事の費用全体が 1 0 0 0 ドルを超えないだろうことを示している。コンテンツ選択ステップ 1 1 5 の実行中、コンテンツ選択装置 5 5 は集合のための記事の選択を制限して、記事の費用全体がこの値を超えないようにしている。この制限により、したがって、集合における記事のソースの関連表示と記事の合計数との両方に影響を与える。

30

【 0 0 8 9 】

なお、ジョブのためのコンテンツ展開の範囲のすべてがそれに追従して実行されなければならない。たとえば、上記に与えられた例を使用して、コンテンツの選択は、(i) どの 1 つのソースも記事の 5 0 % を超えて構成しないように、様々なソースからの記事で構成しなければならず、また、(i i) 調査者のアクセス費を伴う記事（すなわち、ユーザが購読しておらず、または、公共に無料で利用可能ではない記事）を取得するために必要な支出が 1 0 0 0 ドル以下でなければならない。

【 0 0 9 0 】

なお、また、上記に示されたルールは、可能性のあるコンテンツ展開の範囲のルールの一例でしかない。他のタイプのコンテンツ展開の範囲のルールが、本発明の要旨から逸脱することなく、同様に既定および利用され得る。

40

【 0 0 9 1 】

なお、さらに、本明細書ではコンテンツの費用がドルで示されているが、技術を有する専門家が、本発明の要旨を逸脱することなく、国際通貨および関連する通過の兌換で費用のサポートを付けることもできることが理解される。

【 0 0 9 2 】

様々な問合せのルールが既定されると、コンテンツ選択装置 5 5 は、それに続くテキストマイニング操作に使用される予備的な文献の集合を選択する。エンドユーザ 1 3 がテキストマイニングの前に予備的な文献の集合の質を評価することを可能にするために、コン

50

コンテンツ選択装置 55 は、コンテンツ展開の詳細なメトリックを提供するUIレビュースクリーンを生成する。サンプルのUIレビュースクリーン表示は図7(c)に示され、全体として参照符号341で識別される。

【0093】

例示的なスクリーン表示341では、示されたソースのコンテンツ展開が、コンテンツ展開を視覚化するのに有用であるテーブル、またはリスト343、および円グラフ345として提供されている。見て取れるように、集められたコンテンツの42%が単一のソース(無料のソースであるPubMed)から得られる。さらに、集められたコンテンツの約70%が、両方とも無料のソースである、上位2つのソース(PubMedとPLoS)から得られる。

10

【0094】

このことから、ユーザ13はコンテンツ展開が非常に狭い(すなわち、十分なソースが適切に表示されていない)という結論にすぐに達することができる。この結果は、ユーザ15に(i)ソース数が少ないこと、および、(ii)単一のソースが多くを占めていることを知らせる警告347によりサポートされる。

【0095】

ユーザにより、コンテンツ展開が非常に狭く、この理由は、他のことの中で、経費が非常に制限されているためだと判断される場合がある。その結果、ユーザはより良いコンテンツ展開を得るためにコンテンツの費用を増やすことを決める場合がある。

【0096】

20

ユーザにより、コンテンツ展開が非常に狭く、この理由は、他のことの中で、問合せが非常に幅広く、このため最初の文献のプールをかなり多く得たためであると判断する場合もある。その結果、ユーザが検索パラメータの範囲を狭めることを決める場合がある。

【0097】

本明細書にはソースのコンテンツ展開が示されるが、コンテンツ展開の代替的な属性(たとえば、発行日、タイトル、出所の国、記事の言語、および費用の内訳など)が同様に、レビューのためにユーザ13に提供されることが理解される。この相互作用する直観のプロセスを通して、エンドユーザ13は、最終的に最適なコンテンツ展開が得られるまで(たとえば、最適化されたコンテンツ展開が予め既定された経費の範囲に入るまで)、文献の集団を変更することができる。

30

【0098】

最適化されたコンテンツ展開が得られると、具体化されたスケジュールにしたがって、テキストマイニングプロセッサ19によりテキストマイニング操作の処理ステップが実施される。完了すると、結果として得られた書誌的データが収蔵庫21内に導出データセットとして格納される。この獲得されたデータセットは、次いで、プロジェクトマネージャ15を介する調査の中でエンドユーザ13により、引出し、および試験を行うことが可能である。

【0099】

具体的には、ここで図7(d)を参照すると、全体として参照符号351で識別された、テキストマイニングの結果のリストのサンプルスクリーン表示が示されている。見て取れるように、スクリーン表示351には、テキストマイニングプロジェクトの一部として識別される一連の調査文献353-1~353-5のそれぞれの情報(たとえば、書誌的データ、ユーザのアクセス費用、概要など)が含まれる。さらに、このリストに提供された各文献には、ユーザ13に無料または既定の費用で使用可能であれば、記事のフルテキストにアクセスするリンクが含まれる。この方式で、ユーザ13は効率的に、特定のトピックの適切な調査記事をユーザ既定の費用でアクセスおよびレビューすることが可能であり、このことは、本発明の主要な目的である。

40

【0100】

定期的に、エンドユーザ13は様々なテキストおよびデータのマイニングプロジェクトの状態を、プロジェクトマネージャ15によって提供される適切なユーザインターフェイ

50

スを通してレビューおよびモニターすることができる。具体的には、ここで図7(e)を参照すると、エンドユーザ13によって始められた現在および過去のテキストマイニングプロジェクトのレビューのための、ユーザインターフェイスのサンプルスクリーン表示が示されている。この例示的スクリーン表示は全体として参照符号361で識別される。スクリーン表示361では、システム11の、ログインしたエンドユーザ13に利用可能である、開始されたテキストマイニングジョブのテーブル363が示されている。

【0101】

見て取れるように、エンドユーザ13に関する様々なプロジェクトが、コンテンツ選択インターフェイス311を介してユーザによって事前に提供されたプロジェクト名365および詳細情報367を使用して列挙される。さらに、テーブル363は、各プロジェクトのための作品期日ウインドウ369、およびジョブの状態(すなわち、完了した、開かれた、失敗した、処理中など)をユーザに知らせるステータスウインドウ371を含む。さらに、ワンクリック実行ボタン373をクリックすることで、各ジョブに関する特定の機能を取ることができる。

10

【0102】

上記に示された実施形態は、単に例示であることが意図され、当業者であれば、本発明の要旨を逸脱することなく、多数の変形および変更を加えることができる。そのような変形および変更は、添付の特許請求の範囲に既定される本発明の範囲内にあることが意図される。

【符号の説明】

20

【0103】

- 11 システム
- 13 エンドユーザ
- 15 プロジェクトマネージャ
- 17 文献収納庫
- 19 テキストマイニングプロセッサ
- 21 導出データ収納庫
- 55 コンテンツ選択装置

【図 1】

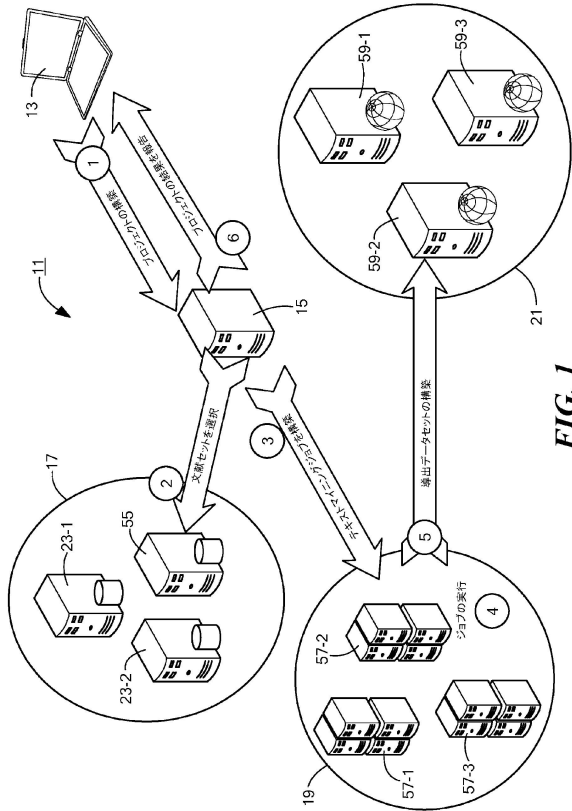


FIG. 1

【図 2】

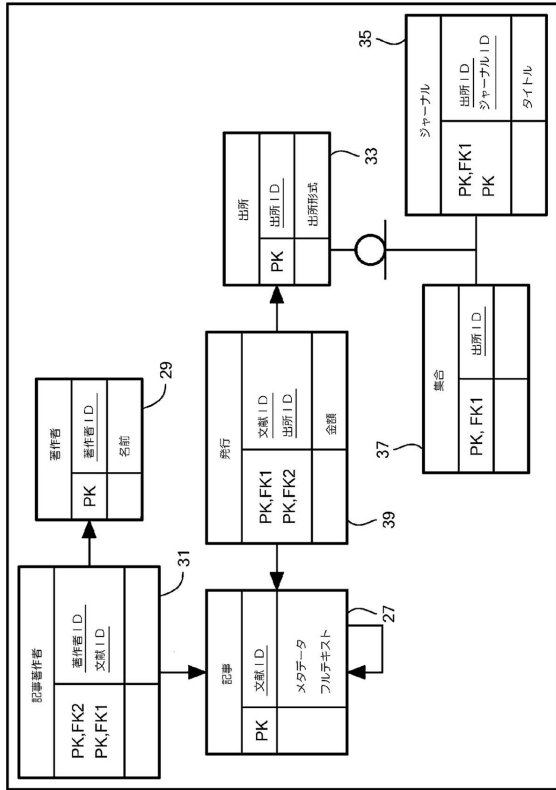


FIG. 2

【図 3】

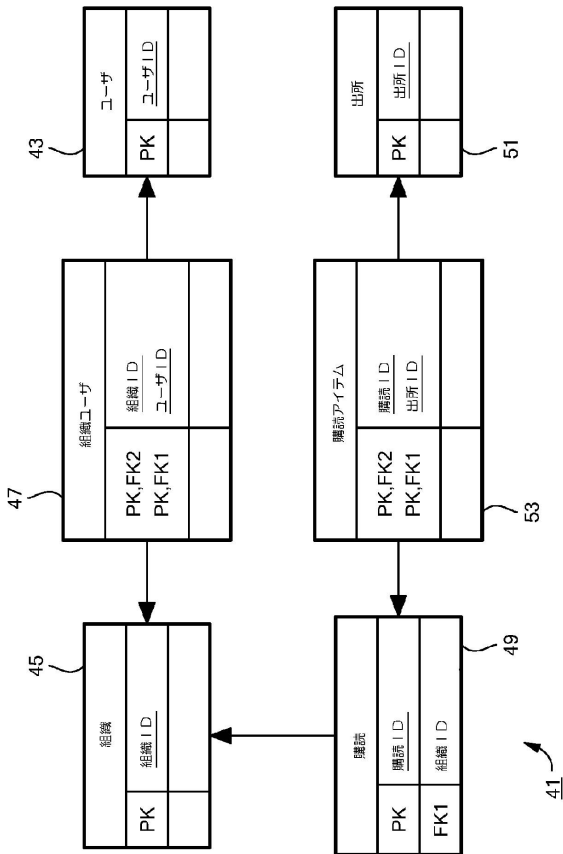


FIG. 3

【図 4】

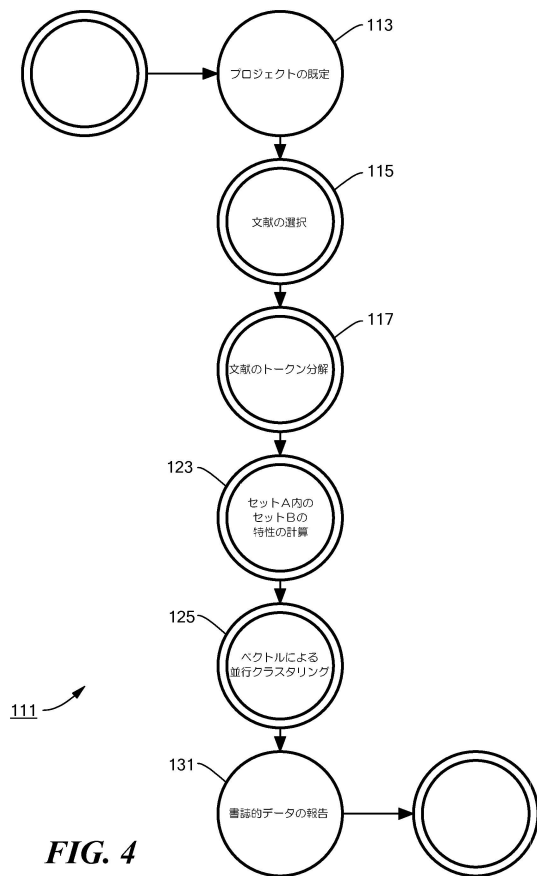


FIG. 4

【図5】

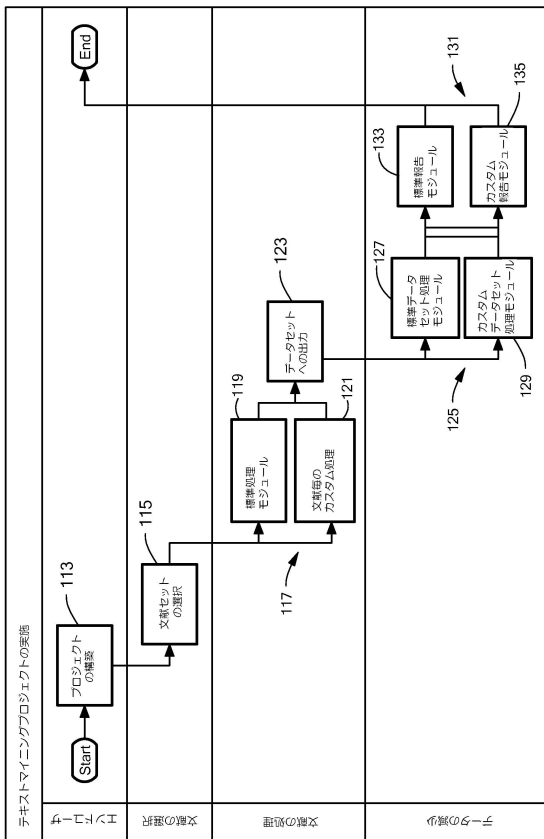


FIG. 5

【図6】

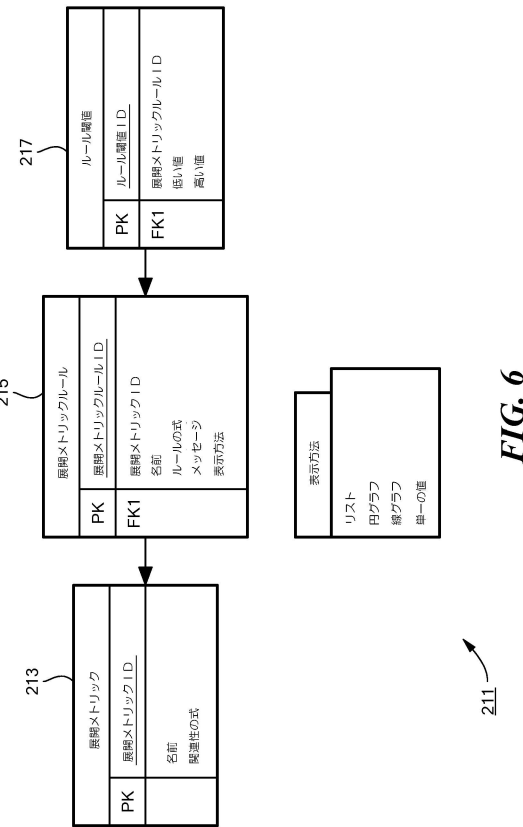


FIG. 6

【図7a】

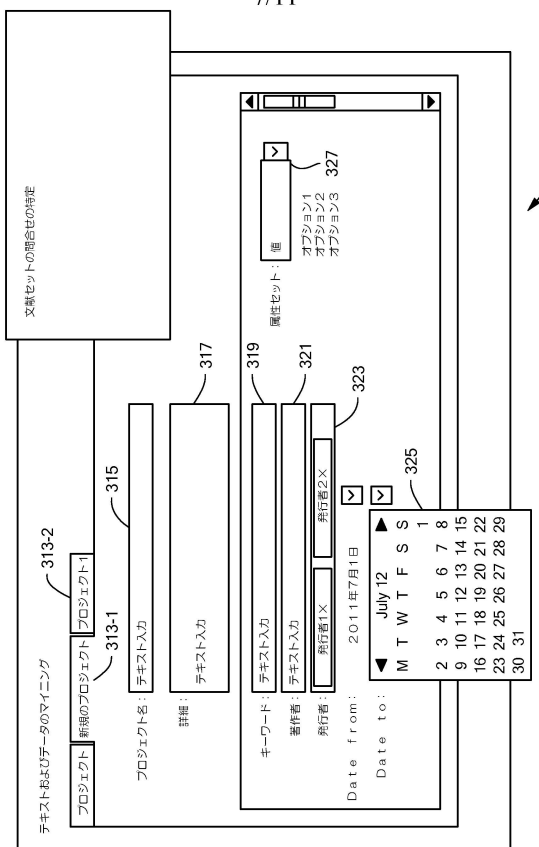


FIG. 7(a)

【図7b】

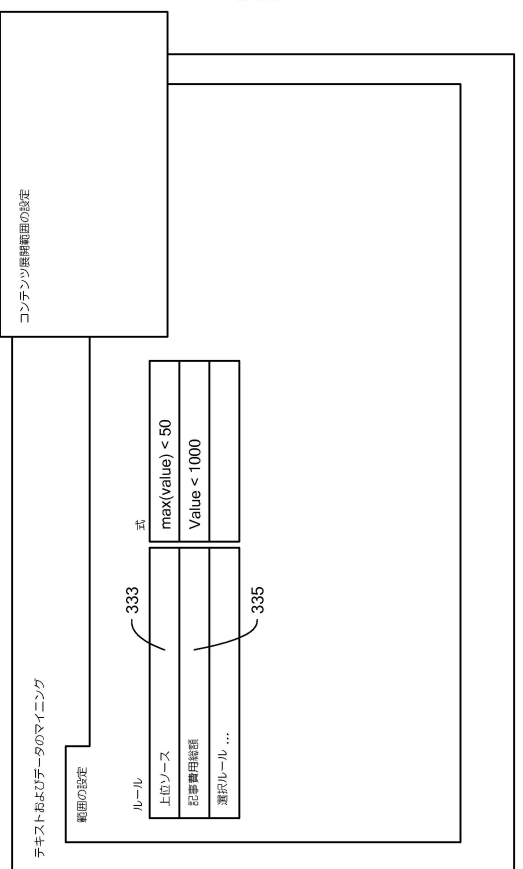


FIG. 7(b)

【 7 c】

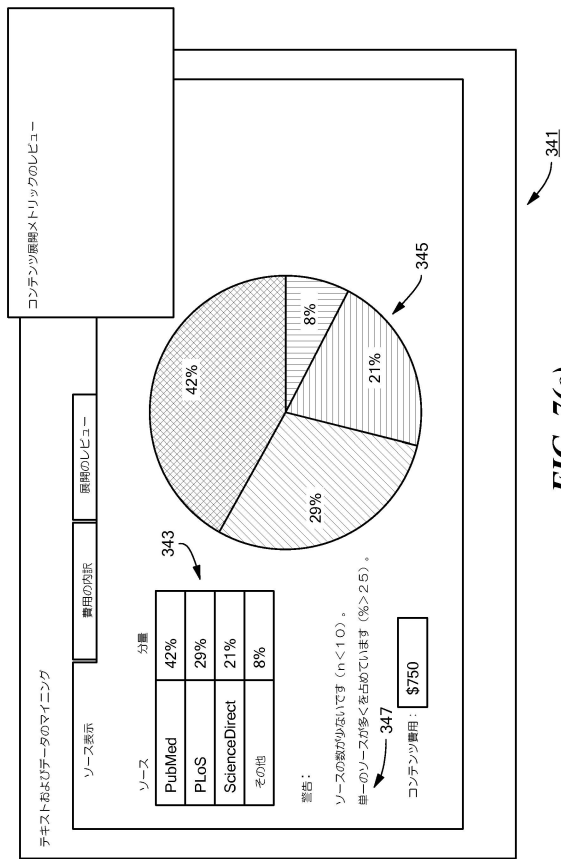


FIG. 7(c)

【 7 d】

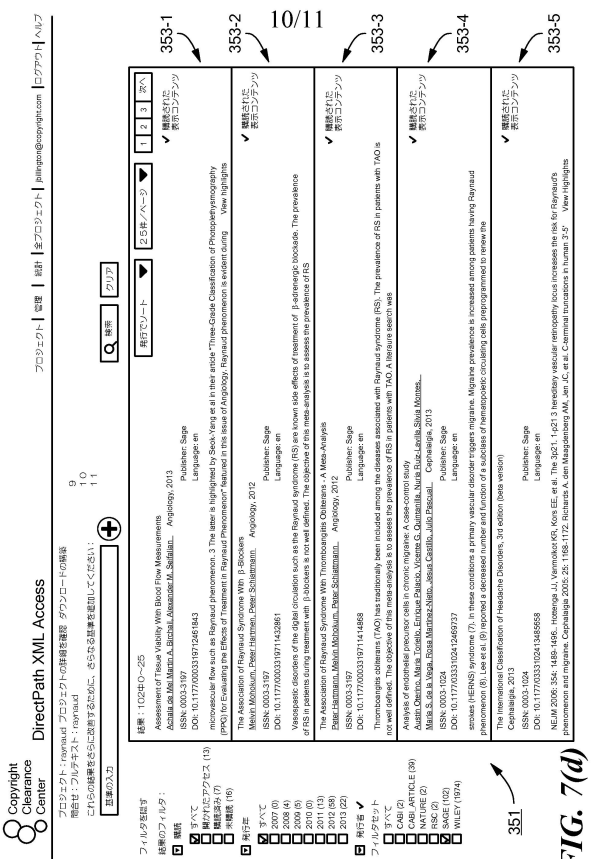


FIG. 7(d)

【 7 e】

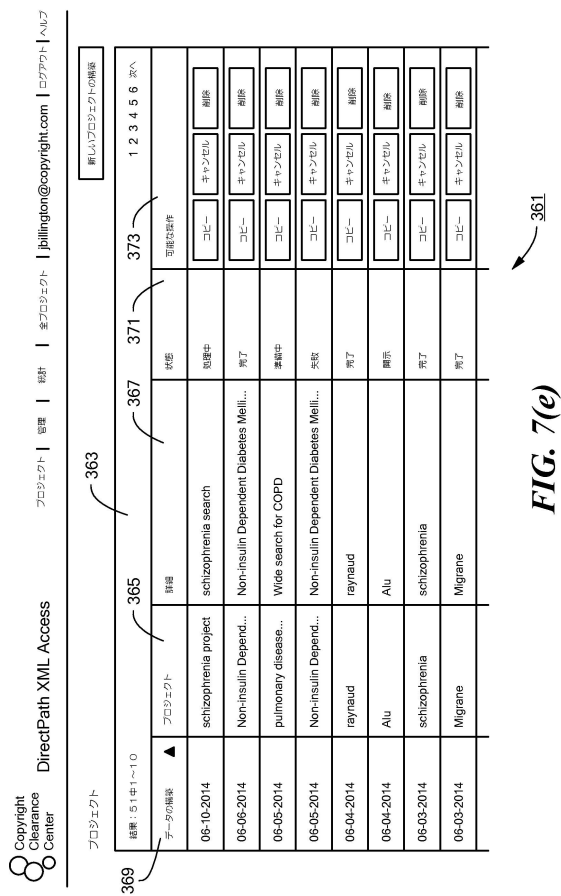


FIG. 7(e)

フロントページの続き

- (72)発明者 クレベ, スコット
アメリカ合衆国 マサチューセッツ州 02420, レキシントン, 8 ハンコック アベニュー
(72)発明者 ビリントン, ジョン
アメリカ合衆国 マサチューセッツ州 01833, ジョージタウン, 42 プロスペクト スト
リート

審査官 松尾 真人

- (56)参考文献 特開2003-216645(JP, A)
特開2009-123139(JP, A)
米国特許出願公開第2012/0221553(US, A1)
特開2008-135057(JP, A)
特表2007-517343(JP, A)
米国特許出願公開第2012/0089642(US, A1)
特表2008-542926(JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06F 17/30