

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号  
特許第5795743号  
(P5795743)

(45) 発行日 平成27年10月14日(2015.10.14)

(24) 登録日 平成27年8月21日(2015.8.21)

(51) Int.Cl.

F I

G O 6 F 17/30 (2006.01)

G O 6 F 17/30 3 5 O C

G O 6 F 17/30 3 4 O Z

G O 6 F 17/30 1 7 O A

請求項の数 5 (全 8 頁)

(21) 出願番号	特願2012-45250 (P2012-45250)	(73) 特許権者	502096543
(22) 出願日	平成24年3月1日(2012.3.1)		パロ・アルト・リサーチ・センター・イン
(65) 公開番号	特開2012-208924 (P2012-208924A)		コーポレーテッド
(43) 公開日	平成24年10月25日(2012.10.25)		P a l o A l t o R e s e a r c h
審査請求日	平成27年2月26日(2015.2.26)		C e n t e r I n c o r p o r a t e d
(31) 優先権主張番号	13/073, 836		アメリカ合衆国、カリフォルニア州 94
(32) 優先日	平成23年3月28日(2011.3.28)		304、パロ・アルト、コヨーテ・ヒル・
(33) 優先権主張国	米国 (US)		ロード 3333
早期審査対象出願		(74) 代理人	100079049
			弁理士 中島 淳
		(74) 代理人	100084995
			弁理士 加藤 和詳
		最終頁に続く	

(54) 【発明の名称】 適応的重み付けを用いた様々な文書間類似度計算方法に基づいた文書比較方法および文書比較システム

(57) 【特許請求の範囲】

【請求項 1】

コンピュータ実行可能な文書比較方法であって、  
コンピュータが、2つの文書に関する少なくとも2つの文書間類似値  $s_i(a, b)$  を受信するステップであって、前記文書間類似値は、様々な文書間類似度計算方法  $i$  によって計算される、前記ステップと、  
2つの文書  $a$  及び  $b$  のそれぞれの個別の文書間類似度計算方法  $i$  の重み  $(w_i, w_i)$  を決定するステップと、  
ここで、重み  $w_i$  は文書  $a$  に対応し、重み  $w_i$  は文書  $b$  に対応し、1つの文書間類似度計算方法の重みを決定するステップは、  
前記重みを初期化するステップと、  
前記重みを、ユーザからのフィードバックに基づいて更新するステップと、を含み、各文書の前記文書間類似度計算方法の前記重みを更新するステップは、ユーザからのフィードバックに基づいて、重みに学習アルゴリズムを適用するステップを含み、  
前記2つの文書に関する前記個別の文書間類似度計算方法  $i$  の組み合わせさせた重みを計算するための重み組み合わせ関数  $f(w_i, w_i)$  を決定するステップと、  
ここで、前記重み組み合わせ関数は、前記文書  $a$  及び  $b$  の間の比較の方向を定め、  
コンピュータが、 $k$  個の文書間類似度計算方法によって計算される前記文書間類似値と前記重み組み合わせ関数とに基づいて、

【数 1】

$$S(a, b) = \sum_{i=1}^k f(\alpha_i, \beta_i) \cdot s_i(a, b)$$

10

として、組み合わせさせた類似値  $S(a, b)$  を生成するステップとを含む、コンピュータ実行可能な文書比較方法。

【請求項 2】

前記文書間類似度計算方法は、テキストに基づいた文書間類似度計算方法、視覚に基づいた文書間類似度計算方法、および、ソーシャルネットワークに基づいた文書間類似度計算方法のうちの 1 つ以上を含む、請求項 1 に記載の方法。

【請求項 3】

各文書の前記文書間類似度計算方法の前記重みは、文書タイプ、文書構造、および、それらの文書の重みのうちの少なくとも 1 つに基づいて初期化される、請求項 1 に記載の方法

20

【請求項 4】

前記重み組み合わせ関数は、前記重みの、平均値、最小値、または、最大値を計算する実数値関数である、請求項 1 に記載の方法。

【請求項 5】

文書間類似度レベルを推定するシステムであって、

文書間類似度レベルを推定するための、少なくとも 1 つのメモリに結合された少なくとも 1 つのハードウェア・プロセッサと、

2 つの文書に関する少なくとも 2 つの文書間類似値  $s_i(a, b)$  を受信するように構成された受信機構であって、前記文書間類似値は、様々な文書間類似度計算方法  $i$  によって計算される、前記受信機構と、

30

2 つの文書  $a$  及び  $b$  のそれぞれの個別の文書間類似度計算方法  $i$  の重み  $(\alpha_i, \beta_i)$  と、

前記 2 つの文書に関する前記個別の文書間類似度計算方法  $i$  の組み合わせさせた重みを計算するための重み組み合わせ関数  $f(\alpha_i, \beta_i)$  と、

を決定するように構成された決定機構であって、

ここで、重み  $\alpha_i$  は文書  $a$  に対応し、重み  $\beta_i$  は文書  $b$  に対応し、1 つの文書間類似度計算方法の重みの決定は、

前記重みを初期化することと、

前記重みを、ユーザからのフィードバックに基づいて更新することと、を含み、各文書の前記文書間類似度計算方法の前記重みを更新することは、ユーザからのフィードバックに基づいて、重みに学習アルゴリズムを適用することを含み、

40

前記重み組み合わせ関数は、前記文書  $a$  及び  $b$  の間の比較の方向を定める、決定機構と、

$k$  個の文書間類似度計算方法によって計算される前記文書間類似値と前記重み組み合わせ関数とに基づいて、

【数 2】

$$S(a, b) = \sum_{i=1}^k f(\alpha_i, \beta_i) \cdot s_i(a, b)$$

10

として、組み合わせさせた類似値  $S(a, b)$  を生成するように構成された、値生成機構とを含む、文書間類似度レベルを推定するシステム。

【発明の詳細な説明】

【背景技術】

【0001】

本開示は、一般的に文書類似度の分析に関する。より詳細には、本開示は、様々な文書間類似度計算方法に基づいた文書比較に関する。

【発明の概要】

【0002】

20

本発明の一実施形態は、適応的重み付けを用いた様々な文書間類似度計算方法に基づいた文書比較システムを提示する。動作中に、システムは、2つの文書に関する少なくとも2つの文書間類似値を受信する。文書間類似値は、様々な文書間類似度計算方法によって計算される。次に、システムは、2つの文書それぞれの個別の文書間類似度計算方法の重みと、2つの文書に関する個別の文書間類似度計算方法の組み合わせさせた重みを計算するための重み組み合わせ関数とを決定する。次に、システムは、文書間類似値と重み組み合わせ関数とに基づいて組み合わせさせた類似値を生成する。

【0003】

本実施形態における一変形形態では、文書間類似度計算方法は、1つまたは複数の、テキストに基づいた、視覚に基づいた、使用に基づいた、および、ソーシャルネットワーク

30

に基づいた文書間類似度計算方法を含む。

【0004】

本実施形態における一変形形態では、文書間類似度計算方法の重みを決定している間、システムは、重みを初期化し、重みをユーザからのフィードバックに基づいて更新する。

【0005】

さらなる一変形形態では、各文書の文書間類似度計算方法の重みは、文書タイプ、文書場所、文書構造、文書使用、および、それらの文書の重みのうちの少なくとも1つに基づいて初期化される。

【0006】

さらなる一変形形態では、各文書の文書間類似度計算方法の重みは、ユーザからのフィードバックに基づいて重みに学習アルゴリズムを適用することによって更新される。

40

【0007】

本実施形態における一変形形態では、重み組み合わせ関数は、重みの、平均値、最小値、または、最大値を計算する実数値関数である。

【図面の簡単な説明】

【0008】

【図1】本発明の一実施形態による類似度組み合わせシステムを示す図である。

【図2】本発明の一実施形態による文書比較プロセスを示すフローチャートである。

【図3】本発明の一実施形態による一文書の文書間類似度計算方法の重みを決定するプロセスを示すフローチャートである。

50

【図 4】本発明の一実施形態による文書比較用の例示的なコンピュータシステムである。

【発明を実施するための形態】

【0009】

本発明の実施形態は、様々な文書間類似度計算方法に基づいた文書比較問題を解決する。動作中に、システムは、2つの文書に関する少なくとも2つの文書間類似値であって様々な文書間類似度計算方法によって計算される文書間類似値を受信する。次に、システムは、2つの文書それぞれの個別の文書間類似度計算方法の重みと、2つの文書に関する個別の文書間類似度計算方法の組み合わせた重みを計算するための重み組み合わせ関数とを決定する。システムは、続いて、文書間類似値と重み組み合わせ関数とに基づいて組み合わせた類似値を生成する。

10

【0010】

メッセージ間または会話間の類似度から導き出すために、計算方法の中には、メッセージ中の、意味のある語または「エンティティ」の出現を比較するものもある。他の方法は、文書使用、または、文書が移動するときに行われる演算の順序を検出することによって、文書類似度を推定する。画像ファイルおよびプレゼンテーションスライドが比較される時、多くの場合、視覚類似度が発見される。幅の広い人気を得たソーシャルネットワークに関しては、ソーシャルネットワークに基づいた方法は、文書間のソーシャル接続に基づいて文書類似度を決定する。

【0011】

本発明の実施形態は、様々な文書タイプ、ユーザの行動様式、および、ユーザの好みを説明しながら、これらの独立した文書間類似度計算方法を効果的に組み合わせるための方法を提示する。例えば、文書 a と文書 b との類似度を計算するとき、k 類似度結果  $s_i(a, b)$  (ここで、 $i = 1 \cdots k$ ) を取得する多数の k 計算方法が展開される。k 結果を組み合わせるために、重み  $\alpha_i$  が文書 a の文書間類似度計算方法 i に割り当てられ、もう1つの重み  $\beta_i$  が文書 b の文書間類似度計算方法 i に割り当てられる。次に、文書間類似度計算方法 i の組み合わせた重みを計算する重み組み合わせ関数  $f(\alpha_i, \beta_i)$  が決定される。組み合わせた類似値  $S(a, b)$  を、

20

【数 1】

$$S(a, b) = \sum_{i=1}^k f(\alpha_i, \beta_i) \cdot s_i(a, b) \quad (1)$$

30

のように計算できる。

【0012】

図 1 は、本発明の一実施形態にかかる類似度組み合わせシステムを示す図である。類似度組み合わせシステム 100 は、組み合わせた類似度計算器 102 と、多数の入力 104、110 と、出力 112 とを含む。動作中に、組み合わせた類似度計算器 102 は、類似度結果  $s_i(a, b)$  の入力 104 と、重み  $\alpha_i$  の入力 106 と、重み  $\beta_i$  の入力 108 と、重み組み合わせ関数  $f(\alpha_i, \beta_i)$  の入力 110 とを受信する。組み合わせた類似度計算器 102 は、次に、入力に基づいて組み合わせた類似値  $S(a, b)$  の出力 112 を計算する。

40

【0013】

文書間類似度計算方法 i の組み合わせた重みを計算するための重み組み合わせ関数  $f(\alpha_i, \beta_i)$  は、重みのパラメータの、平均値、最小値、または、最大値を計算する実数値関数といった、様々な形態をとっていてもよい。例えば、 $f(\alpha_i, \beta_i)$  は、 $\alpha_i$  と  $\beta_i$  との線形結合、

【数 2】

$$f(\alpha_i, \beta_i) = x \cdot \alpha_i + y \cdot \beta_i \quad (2)$$

50

ここで  $x + y = 1$  として規定されてもよい。 $x = y$  であれば、 $f(i, i)$  は、 $i$  と  $i$  との平均値に等しい。さらに、文書  $a$  と文書  $b$  との類似度が対称的ではない、すなわち  $s_i(a, b) \neq s_i(b, a)$  ならば、関数  $f(i, i)$  は、設定  $x = y$  によってこのことを説明することができる。具体的には、 $x > y$  の場合、組み合わせさせた重み  $f(i, i)$  は、文書  $a$  から文書  $b$  への比較方向を強調する。

【0014】

図2は、本発明の一実施形態による文書比較プロセスを示すフローチャートである。動作中に、システムは、様々な文書間類似度計算方法によって計算される多数の類似値を受信する(演算202)。次に、システムは、各文書間類似度計算方法の重みを決定する(演算204)。次に、システムは、重み組み合わせ関数を決定する(演算206)。続いて、システムは、組み合わせさせた類似度を生成する(演算208)。

10

【0015】

文書間類似度計算方法の重みを決定するとき、文書比較システムは、重み値を初期化して、重み値をユーザからのフィードバックに基づいて更新する必要がある。図3は、本発明の一実施形態にかかる文書の文書間類似度計算方法の重みを決定するプロセスを示すフローチャートである。動作中に、システムは、文書タイプ、文書場所、および、文書構造に基づいて、重み値を初期化し(演算302)、重みをユーザからのフィードバックに基づいて更新する(演算304)。

【0016】

各文書は、様々な計算方法に関する文書間類似度方法の重みのベクトルを有している。このベクトルの初期値を、文書タイプ、文書場所、および、文書構造から導き出すことができる。重みベクトルをまた、それらの文書の重み、それらの文書の場所、または、文書使用の周波数に基づいて初期化できる。

20

【0017】

しかし、初めに割り当てられた重み値は、正確ではない場合があり、その場合は使用中に調整される必要がある。極めて重要なのは、ユーザが、類似度組み合わせ結果を順位付けし、これらの結果の不一致を指摘することによって、フィードバックを提示することである。ユーザからのフィードバックに基づいて、重みを更新でき、各類似度計算方法に関して絶えず精密化できる。このために、類似値を入力と捉え、重みを出力と捉え、ユーザからのフィードバックをグラウンドトゥースと捉える、機械学習アルゴリズムを適用できる。一実施形態では、システムは、AdaBoostといったブースティングアルゴリズムを適用して、重みを確認および更新する。AdaBoostは、不一致の結果として選り抜かれたこれらの実現値を支持して、重みを適応できるように微調整する。従って、より精密化された重みは、不一致を正すことができる。

30

【0018】

図4は、本発明の一実施形態による文書比較のための例示的なコンピュータシステムを示す。一実施形態では、コンピュータ通信システム400は、プロセッサ402と、メモリ404と、記憶装置406とを含む。記憶装置406は、文書比較アプリケーション408、および、アプリケーション410、412といった他のアプリケーションを記憶する。動作中に、文書比較アプリケーション408を、記憶装置406からメモリ404に読み込み、次に、プロセッサ402によって実行する。プログラムを実行している間、プロセッサ402は、前述の関数を行う。コンピュータ通信システム400は、任意の表示装置414と、キーボード416と、ポインティングデバイス418とに連結している。

40

【0019】

本詳細な説明において説明したデータ構造およびコードは、通常、コンピュータ読取り可能記憶媒体に格納されている。記憶媒体は、コンピュータシステムによって用いられるコードおよび/またはデータを格納できる任意の装置または媒体であってもよい。コンピュータ読取り可能記憶媒体は、揮発性メモリ、不揮発性メモリ、ディスクドライブなどの磁気記憶装置および光学記憶装置、磁気テープ、CD(コンパクトディスク)、DVD(

50

デジタル多用途ディスクまたはデジタルビデオディスク)、または、現在知られているか、または、今後開発されるコンピュータ読取り可能媒体を格納できる他の媒体を含むが、これらに限定されるものではない。

#### 【0020】

本詳細な説明部分において説明した方法およびプロセスを、上述したようなコンピュータ読取り可能記憶媒体に格納されうるコードおよび/またはデータとして実施できる。コンピュータシステムが、コンピュータ読取り可能記憶媒体に格納されたコードおよび/またはデータを読み取り、実行するとき、コンピュータシステムは、データ構造およびコードとして実施され、コンピュータ読取り可能記憶媒体内に格納された方法およびプロセスを行う。

10

#### 【0021】

さらに、本発明において説明した方法およびプロセスを、ハードウェアモジュールまたはハードウェア機器に含むことができる。これらのモジュールまたは機器は、特定用途向けIC(AASIC)チップ、書替え可能ゲートアレイ(FPGA)、特定のソフトウェアモジュールまたは特定の時間においてコードの一部を実行する専用プロセッサまたは共用プロセッサ、および/または、現在知られているまたは今後開発される他のプログラマブルロジックデバイスを含んでいてもよいが、それらに制限されるものではない。ハードウェアモジュールまたはハードウェア機器を起動するとき、それらは、それらに含まれる方法およびプロセスを行う。

#### 【0022】

20

様々な複数の実施形態に関する上述の説明は、図解および説明するためにのみ示されたものである。これらの説明は、完全であること、あるいは、開示した形態に本発明を限定することを意図したものではない。従って、多くの変更形態および変形形態が、当業者には明らかであろう。さらに、上述の開示は、本発明を限定することを意図したものではない。

【図1】

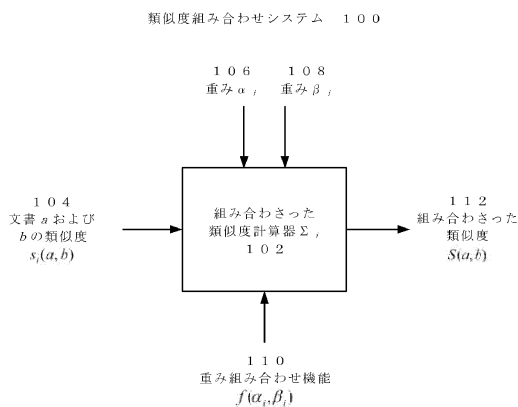


図 1

【図2】

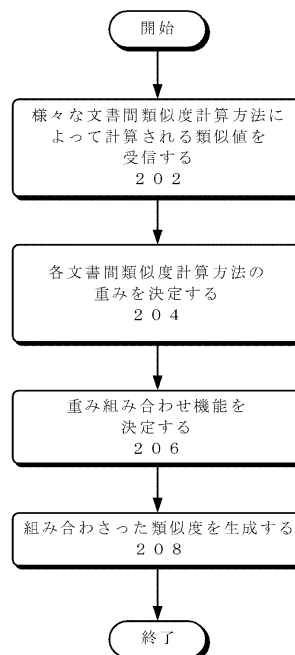


図 2

【図 3】

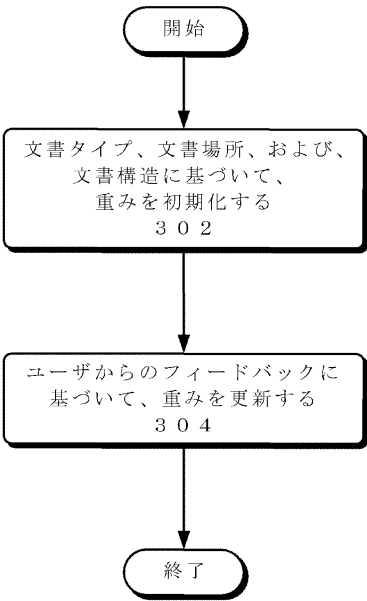


図 3

【図 4】

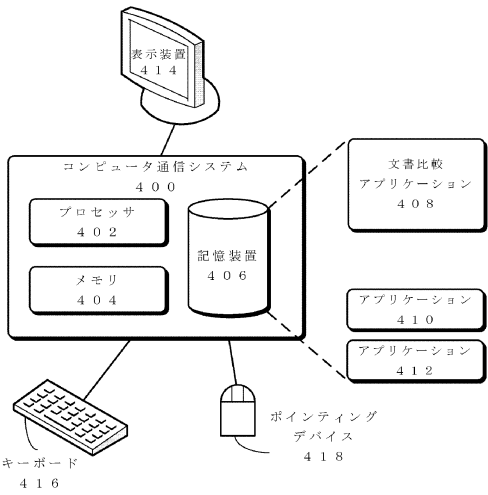


図 4

---

フロントページの続き

(72)発明者 オリヴァー・ブルディスカ  
アメリカ合衆国 カリフォルニア州 9 4 3 0 4 マウンテン・ビュー オルテガ・アヴェニュー  
5 6 5 ナンバー 1 1

(72)発明者 ペトロ・イザレフ  
アメリカ合衆国 カリフォルニア州 9 4 3 0 9 パロ・アルト シェリダン・アヴェニュー 4  
1 0 アpartment 3 4 0

審査官 田中 秀樹

(56)参考文献 特開 2 0 0 0 - 2 0 0 2 8 5 ( J P , A )  
特開平 0 9 - 2 1 2 5 0 9 ( J P , A )  
特開 2 0 0 1 - 1 5 5 0 2 7 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)  
G 0 6 F 1 7 / 3 0