

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7464193号
(P7464193)

(45)発行日 令和6年4月9日(2024.4.9)

(24)登録日 令和6年4月1日(2024.4.1)

(51)国際特許分類 F I
G 0 6 F 16/903(2019.01) G 0 6 F 16/903

請求項の数 10 (全27頁)

(21)出願番号	特願2023-520672(P2023-520672)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和3年5月13日(2021.5.13)	(74)代理人	100103090 弁理士 岩壁 冬樹
(86)国際出願番号	PCT/JP2021/018169	(74)代理人	100124501 弁理士 塩川 誠人
(87)国際公開番号	WO2022/239174	(72)発明者	大野 善之 東京都港区芝五丁目7番1号 日本電気 株式会社内
(87)国際公開日	令和4年11月17日(2022.11.17)	審査官	齋藤 貴孝
審査請求日	令和5年9月13日(2023.9.13)		

最終頁に続く

(54)【発明の名称】 類似度導出システムおよび類似度導出方法

(57)【特許請求の範囲】

【請求項1】

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算手段と、

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定手段と、

前記複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出手段とを備える

ことを特徴とする類似度導出システム。

【請求項2】

前記ハッシュ値計算手段は、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第1のハッシュ値を計算する第1のハッシュ値計算手段と、

前記各集合の個々の要素のうち、前記第1のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは1つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む1つの集合である全

10

20

体集合を生成する全体集合生成手段と、

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成手段と、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算手段と、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第1のハッシュ値と、当該第1のハッシュ値に対応する、前記第2のハッシュ値計算手段によって計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定手段とを含む

10

請求項1に記載の類似度導出システム。

【請求項3】

前記ハッシュ値計算手段は、

前記複数の集合から、順次、1つの集合を選択する集合選択手段と、

選択された集合から、順次、1つの要素を選択する要素選択手段と、

選択された要素に前記所定のハッシュ関数を適用することによって前記選択された要素の第1のハッシュ値を計算する第1のハッシュ値計算手段と、

前記選択された要素と、第1のハッシュ値が一致し、かつ、要素自体が一致する要素である一致要素が既に選択されているか否かを判定する判定手段と、

前記一致要素が既に選択されている場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を、前記一致要素の前記各ハッシュ関数に対応するハッシュ値と同一であると定め、

20

前記一致要素が選択されていない場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算手段とを含む

請求項1に記載の類似度導出システム。

【請求項4】

前記所定のハッシュ関数は、前記複数のハッシュ関数の中で最も値域が広いハッシュ関数である

請求項1から請求項3のうちのいずれか1項に記載の類似度導出システム。

30

【請求項5】

前記類似度導出手段は、

個々のハッシュ関数に対応する最小ハッシュ値同士が一致している数を、ペアをなす2つの集合の類似度として定める

請求項1から請求項4のうちのいずれか1項に記載の類似度導出システム。

【請求項6】

コンピュータが、

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算処理、

40

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理、および、

前記複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理を実行する

ことを特徴とする類似度導出方法。

【請求項7】

前記コンピュータが、

50

前記ハッシュ値計算処理で、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第1のハッシュ値を計算する第1のハッシュ値計算処理、

前記各集合の個々の要素のうち、前記第1のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは1つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む1つの集合である全体集合を生成する全体集合生成処理、

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成処理、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算処理、および、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第1のハッシュ値と、当該第1のハッシュ値に対応する、前記第2のハッシュ値計算処理で計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定処理を実行する

請求項6に記載の類似度導出方法。

【請求項8】

前記コンピュータが、

前記ハッシュ値計算処理で、

前記複数の集合から、順次、1つの集合を選択する集合選択処理、

選択された集合から、順次、1つの要素を選択する要素選択処理、

選択された要素に前記所定のハッシュ関数を適用することによって前記選択された要素の第1のハッシュ値を計算する第1のハッシュ値計算処理、

前記選択された要素と、第1のハッシュ値が一致し、かつ、要素自体が一致する要素である一致要素が既に選択されているか否かを判定する判定処理、および、

前記一致要素が既に選択されている場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を、前記一致要素の前記各ハッシュ関数に対応するハッシュ値と同一であると定め、

前記一致要素が選択されていない場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算処理を実行する

請求項6に記載の類似度導出方法。

【請求項9】

コンピュータに、

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算処理、

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理、および、

前記複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理

を実行させるための類似度導出プログラム。

【請求項10】

前記コンピュータに、

前記ハッシュ値計算処理で、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第1

10

20

30

40

50

のハッシュ値を計算する第1のハッシュ値計算処理、

前記各集合の個々の要素のうち、前記第1のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは1つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む1つの集合である全体集合を生成する全体集合生成処理、

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成処理、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算処理、および、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第1のハッシュ値と、当該第1のハッシュ値に対応する、前記第2のハッシュ値計算処理で計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定処理を実行させる

請求項9に記載の類似度導出プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、集合同士の類似度を導出する類似度導出システム、類似度導出方法、および、類似度導出プログラムに関する。

【背景技術】

【0002】

集合の各要素に同一のハッシュ関数を適用することによって、要素毎にハッシュ値が得られる。そのハッシュ値の最小値をMinHashと称する場合がある。本明細書では、同一のハッシュ関数に基づいて要素毎に得られるハッシュ値の最小値を最小ハッシュ値と記す。前述の集合の各要素に複数のハッシュ関数を適用すれば、ハッシュ関数と同数の最小ハッシュ値が得られる。

【0003】

また、2つの集合に、共通の複数のハッシュ関数を適用し、集合毎に複数の最小ハッシュ値を求め、その複数の最小ハッシュ値に基づいて、2つの集合の類似度を求めることが考えられる。

【0004】

また、特許文献1には、最小ハッシュ値を用いて類似テキストを探索する方法が記載されている。

【先行技術文献】

【特許文献】

【0005】

【文献】特開2020-4107号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

前述の例のように集合の類似度を求めるためには、個々の集合の個々の要素毎に、複数のハッシュ関数を適用することによって複数のハッシュ値を求めなければならない。これは、集合毎に、複数の最小ハッシュ値を求める必要があるためである。ここで、集合の個数を n 個とする。また、ここでは、説明を簡単にするために、各集合の要素数が k 個で共通であるとする。また、ハッシュ関数の数を m 個とする。この場合、ハッシュ値の計算を、 $n \cdot k \cdot m$ 回行う必要がある。

【0007】

しかし、ハッシュ値の計算量は少ない方が好ましい。

【0008】

10

20

30

40

50

そこで、本発明は、集合の類似度を導出する際に、ハッシュ値の計算量を低減させることができる類似度導出システム、類似度導出方法、および、類似度導出プログラムを提供することを目的とする。

【課題を解決するための手段】

【0009】

本発明による類似度導出システムは、複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、その複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求めるハッシュ値計算手段と、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定手段と、複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出手段とを備えることを特徴とする。

10

【0010】

本発明による類似度導出方法は、コンピュータが、複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、その複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求めるハッシュ値計算処理、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理、および、複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理を実行することを特徴とする。

20

【0011】

本発明による類似度導出プログラムは、コンピュータに、複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、その複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求めるハッシュ値計算処理、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理、および、複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理を実行させる。

30

【発明の効果】

【0012】

本発明によれば、集合の類似度を導出する際に、ハッシュ値の計算量を低減させることができる。

40

【図面の簡単な説明】

【0013】

【図1】本発明の第1の実施形態の類似度導出システムの例を示すブロック図である。

【図2】第1の実施形態の処理経過の例を示すフローチャートである。

【図3】本発明の第2の実施形態の類似度導出システムの例を示すブロック図である。

【図4】各集合A、Bの個々の要素毎に得られた第1のハッシュ値の例を示す模式図である。

【図5】図4に示す集合A、Bから得られる全体集合を示す模式図である。

【図6】インデックス集合の例を示す模式図である。

50

【図 7】図 5 に示す全体集合に属する各要素に対して、第 2 のハッシュ値計算部 5 4 が計算した各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の例を示す模式図である。

【図 8】各集合の個々の要素に対して決定された、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値の例を示す模式図である。

【図 9】集合 A、集合 B に関してそれぞれ、ハッシュ関数毎に特定された最小ハッシュ値を示す模式図である。

【図 10】第 2 の実施形態の処理経過の例を示すフローチャートである。

【図 11】ステップ S 3 の動作の具体例の一例を示すフローチャートである。

【図 12】本発明の第 3 の実施形態の類似度導出システムの例を示すブロック図である。

【図 13】第 3 の実施形態の処理経過の例を示すフローチャートである。

10

【図 14】第 3 の実施形態の処理経過の例を示すフローチャートである。

【図 15】複数の集合のうちの 1 つの集合の例を示す模式図である。

【図 16】最初にステップ S 7 3 を実行した後に得られたハッシュ値の例を示す模式図である。

【図 17】最初にステップ S 7 6 を実行した後に得られたハッシュ値の例を示す模式図である。

【図 18】ステップ S 7 7 で未選択の要素がないと判定されるまでに得られたハッシュ値の例を示す模式図である。

【図 19】複数の集合のうちの 1 つの集合の例を示す模式図である。

【図 20】集合 B の要素 "mountain" が選択され、ステップ S 7 3 を実行した後に得られたハッシュ値の例を示す模式図である。

20

【図 21】集合 B の要素 "mountain" を選択したときにおけるステップ S 7 5 の実行後に得られているハッシュ値の例を示す模式図である。

【図 22】各集合の各要素に対して求められた複数のハッシュ値の例を示す模式図である。

【図 23】本発明の実施形態の類似度導出システム 1 に係るコンピュータの構成例を示す概略ブロック図である。

【発明を実施するための形態】

【0014】

以下、本発明の実施形態を図面を参照して説明する。以下に示す各実施形態では、複数の集合が、予め、各実施形態の類似度導出システムに入力されているものとする。

30

【0015】

実施形態 1 .

第 1 の実施形態は、本発明の概要を示す実施形態である。より具体的な事項については、後述の第 2 の実施形態、第 3 の実施形態で説明する。

【0016】

図 1 は、本発明の第 1 の実施形態の類似度導出システムの例を示すブロック図である。類似度導出システム 1 は、ハッシュ値計算部 2 と、最小ハッシュ値特定部 3 と、類似度導出部 4 とを備える。

【0017】

ハッシュ値計算部 2 は、複数の集合に含まれる各集合の個々の要素に対して、複数のハッシュ関数を適用して得られる複数のハッシュ値を求める。ただし、このとき、ハッシュ値計算部 2 は、その複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求める。

40

【0018】

複数のハッシュ関数の中で値域が最も広いハッシュ関数を所定のハッシュ関数として定めることが好ましい。すなわち、所定のハッシュ関数は、複数のハッシュ関数の中で最も値域が広いハッシュ関数であることが好ましい。

【0019】

50

最小ハッシュ値特定部 3 は、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する。例えば、1つの集合 A に着目したとする。また、あるハッシュ関数に着目した場合、そのハッシュ関数に対応する、集合 A の要素数分のハッシュ値が得られている。最小ハッシュ値特定部 3 は、そのハッシュ関数に対応する最小ハッシュ値として、それらのハッシュ値の最小値を特定する。最小ハッシュ値特定部 3 は、他の各ハッシュ関数に対してもそれぞれ、同様に、最小ハッシュ値を特定する。最小ハッシュ値特定部 3 は、この処理を、複数の集合に含まれるそれぞれの集合に対して行う。この結果、集合毎に、個々のハッシュ関数に対応する最小ハッシュ値が定まる。

【 0 0 2 0 】

10

類似度導出部 4 は、複数の集合から得られる 1 組以上の集合のペアに関して、ペアをなす 2 つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する。例えば、類似度導出部 4 は、複数の集合から得られる全てのペア（集合のペア）に関して、ペアをなす 2 つの集合の類似度を算出してもよい。

【 0 0 2 1 】

また、類似度導出部 4 は、例えば、個々のハッシュ値に対応する最小ハッシュ値同士が一致している数を、ペアをなす 2 つの集合の類似度として定めてもよい。

【 0 0 2 2 】

ハッシュ値計算部 2、最小ハッシュ値特定部 3、および、類似度導出部 4 は、例えば、類似度導出プログラムに従って動作するコンピュータの CPU (Central Processing Unit) によって実現される。例えば、CPU が、コンピュータのプログラム記憶装置等のプログラム記録媒体から類似度導出プログラムを読み込み、その類似度導出プログラムに従って、ハッシュ値計算部 2、最小ハッシュ値特定部 3、および、類似度導出部 4 として動作すればよい。

20

【 0 0 2 3 】

図 2 は、第 1 の実施形態の処理経過の例を示すフローチャートである。なお、既に説明した事項については、適宜、説明を省略する。

【 0 0 2 4 】

まず、ハッシュ値計算部 2 が、所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求める（ステップ S 1）。

30

【 0 0 2 5 】

次に、最小ハッシュ値特定部 3 が、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に最小ハッシュ値を特定する（ステップ S 2）。

【 0 0 2 6 】

次に、類似度導出部 4 が、複数の集合から得られる 1 組以上の集合のペアに関して、ペアをなす 2 つの集合の類似度を導出する（ステップ S 3）。

【 0 0 2 7 】

本実施形態によれば、ハッシュ値計算部 2 は、複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ値を求める際に、所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除する。従って、その複数の要素に対して、同一のハッシュ関数の計算が繰り返し行われることはない。よって、集合の類似度を導出する際に、ハッシュ値の計算量を低減させることができる。

40

【 0 0 2 8 】

実施形態 2 .

第 2 の実施形態は、第 1 の実施形態をより具体的に示した実施形態である。図 3 は、本発明の第 2 の実施形態の類似度導出システムの例を示すブロック図である。本実施形態の類似度導出システム 1 も、ハッシュ値計算部 2 と、最小ハッシュ値特定部 3 と、類似度導

50

出部 4 とを備える。最小ハッシュ値特定部 3 および類似度導出部 4 は、第 1 の実施形態における最小ハッシュ値特定部 3 および類似度導出部 4 と同様である。

【 0 0 2 9 】

また、ハッシュ値計算部 2 は、第 1 のハッシュ値計算部 5 1 と、全体集合生成部 5 2 と、インデックス情報生成部 5 3 と、第 2 のハッシュ値計算部 5 4 と、ハッシュ値決定部 5 5 とを備える。

【 0 0 3 0 】

第 1 のハッシュ値計算部 5 1 は、複数の集合に含まれる各集合の個々の要素に対して、複数のハッシュ関数のうちの所定のハッシュ関数を適用することによって、ハッシュ値を計算する。この所定のハッシュ関数によって得られるハッシュ値を第 1 のハッシュ値と記す。したがって、各集合の個々の要素毎に第 1 のハッシュ値が計算される。

10

【 0 0 3 1 】

以下では、説明を簡単にするために 2 つの集合 A , B に着目して説明するが、集合の数は 3 つ以上であってもよい。また、複数のハッシュ関数の数は、m 個であり、個々のハッシュ関数を h_1, h_2, \dots, h_m と記すこととする。そして、 h_1 が、上記の所定のハッシュ関数であるものとする。

【 0 0 3 2 】

なお、複数のハッシュ関数の中から 1 つハッシュ関数を、予め所定の関数として定めておけばよい。このとき、複数のハッシュ関数の中で値域が最も広いハッシュ関数を所定のハッシュ関数として定めることが好ましい。すなわち、所定のハッシュ関数は、複数のハッシュ関数の中で最も値域が広いハッシュ関数であることが好ましい。本例では、ハッシュ関数 h_1 の値域が、複数のハッシュ関数 $h_1 \sim h_m$ の値域の中で最も広いものとする。この点は、前述の第 1 の実施形態や後述の第 3 の実施形態においても同様である。

20

【 0 0 3 3 】

図 4 は、各集合 A , B の個々の要素毎に得られた第 1 のハッシュ値の例を示す模式図である。図 4 に示すように、集合 A は 4 つ文字列を要素として持ち、集合 B は、3 つの文字列を要素として持っているものとする。第 1 のハッシュ値計算部 5 1 は、これらの個々の要素に対して、所定のハッシュ関数 h_1 によって、第 1 のハッシュ値を計算する。図 4 では、個々の要素の右側に示した値が、第 1 のハッシュ値である。

【 0 0 3 4 】

全体集合生成部 5 2 は、複数の集合に含まれる各集合の個々の要素のうち、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、その複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む 1 つの集合を生成する。以下、この集合を全体集合と記す。

30

【 0 0 3 5 】

図 5 は、図 4 に示す集合 A , B から得られる全体集合を示す模式図である。図 4 に示すように、各集合の個々の要素のうち、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素は、集合 A に属する " mountain " と、集合 B に属する " mountain " である。この 2 つの要素は、第 1 のハッシュ値が " 6 6 6 " で一致し、かつ、要素自体が " mountain " で一致している。全体集合生成部 5 2 は、この 2 つの要素から 1 つの要素だけを取り出す。このとき、集合 A に属する " mountain " を取りだしてもよく、あるいは、集合 B に属する " mountain " を取りだしてもよい。そして、全体集合生成部 5 2 は、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素 (集合 A に属する " mountain " 、および、集合 B に属する " mountain ") に該当しない各要素 (本例では、" the " , " highest " , " Fuji " , " room " , " view ") を取り出す。そして、全体集合生成部 5 2 は、取り出した各要素を含む 1 つの集合を、全体集合として生成する。本例では、図 5 に示すように、全体集合は、要素として、6 個の文字列 (" the " , " highest " , " mountain " , " Fuji " , " room " , " view ") を含む。全体集合生成部 5 2 は、前述のように集合 A , B から要素を取り出しているため、全体集合に属する要素に重複はなく、また、全体集合は、集合 A , B に属する全ての種類の要素を含んでいる。

40

50

【 0 0 3 6 】

インデックス情報生成部 5 3 は、複数の集合に含まれるどの集合のどの要素が全体集合のどの要素に該当するかを示すインデックス情報を生成する。

【 0 0 3 7 】

図 6 は、インデックス集合の例を示す模式図である。図 6 に例示するインデックス情報を生成する例を説明する。インデックス情報生成部 5 3 は、全体集合に属する各要素に対して識別情報を割り当てる。図 6 に示す例では、全体集合に属する各要素に対して “ 1 ” から “ 6 ” の識別情報が割り当てられている。そして、インデックス情報生成部 5 3 は、複数の集合に含まれる各集合の各要素毎に、その要素に該当する全体集合内の要素の識別情報を定めた情報を、インデックス情報として生成する。図 6 に示す例では、例えば、集合 A に属する要素 “ the ” には、“ 1 ” というインデックス情報が定められている。このことは、集合 A に属する要素 “ the ” が、全体集合内の要素のうちの、識別要素 “ 1 ” が割り当てられた要素に該当することを示している。また、図 6 に示す例では、集合 A に属する要素 “ mountain ” 、および、集合 B に属する要素 “ mountain ” は、いずれも、全体集合内の要素のうちの、識別要素 “ 3 ” が割り当てられた要素に該当することを示している。すなわち、集合 A に属する要素 “ mountain ” 、および、集合 B に属する要素 “ mountain ” は、全体集合内の同一の要素に該当することを示している。

10

【 0 0 3 8 】

第 2 のハッシュ値計算部 5 4 は、全体集合に属する各要素に対して、複数のハッシュ関数（本例では、ハッシュ関数 h_1, h_2, \dots, h_m ）のうちの所定のハッシュ関数（本例では、ハッシュ関数 h_1 ）以外の各ハッシュ関数（すなわち、各ハッシュ関数 h_2, \dots, h_m ）を適用する。そして、第 2 のハッシュ値計算部 5 4 は、全体集合に属する各要素に対して、その各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値をそれぞれ計算する。

20

【 0 0 3 9 】

図 7 は、図 5 に示す全体集合に属する各要素に対して、第 2 のハッシュ値計算部 5 4 が計算した各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の例を示す模式図である。

【 0 0 4 0 】

ハッシュ値決定部 5 5 は、インデックス情報に基づいて、複数の集合に含まれる各集合（本例では、集合 A , B）の個々の要素の第 1 のハッシュ値（図 4 参照）と、その第 1 のハッシュ値に対応する、第 2 のハッシュ値計算部 5 4 によって計算されたハッシュ値（図 7 参照）とを組み合わせることによって、各集合の個々の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値を決定する。

30

【 0 0 4 1 】

本例では、各集合の各要素毎に、その要素に該当する全体集合内の要素の識別情報を定めた情報を、インデックス情報としている（図 6 参照）。また、第 1 のハッシュ値計算部 5 1 によって、各集合の各要素に対応付けて第 1 のハッシュ値が計算されている（図 4 参照）。さらに、第 2 のハッシュ値計算部 5 4 によって、全体集合に属する各要素に対して、所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値が計算されている。従って、ハッシュ値決定部 5 5 は、インデックス情報に基づいて、各集合の個々の要素の第 1 のハッシュ値と、第 2 のハッシュ値計算部 5 4 によって計算された各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を組み合わせることができる。そして、その結果、ハッシュ値決定部 5 5 は、各集合の個々の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値を決定することができる。

40

【 0 0 4 2 】

図 8 は、各集合（本例では、集合 A , B）の個々の要素に対して決定された、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値の例を示す模式図である。

50

【 0 0 4 3 】

このように、第 1 のハッシュ値計算部 5 1、全体集合生成部 5 2、インデックス情報生成部 5 3、第 2 のハッシュ値計算部 5 4、および、ハッシュ値決定部 5 5 の動作によって、各集合の個々の要素に対して、複数のハッシュ関数を適用して得られる複数のハッシュ値を求める。また、このとき、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素（本例では、集合 A に属する " mountain "、および、集合 B に属する " mountain "）については、所定のハッシュ関数 h_1 以外のハッシュ関数 h_2, \dots, h_m の計算を重複して行っていない。具体的には、" mountain " に対する、ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の計算は、第 2 のハッシュ値計算部 5 4 がそれぞれ 1 回行っているのみである。

10

【 0 0 4 4 】

最小ハッシュ値特定部 3 は、複数の集合に含まれるそれぞれの集合（本例では、集合 A、集合 B）に関して、個々のハッシュ関数毎に最小ハッシュ値を特定する。図 9 は、集合 A、集合 B に関してそれぞれ、ハッシュ関数毎に特定された最小ハッシュ値を示す模式図である。例えば、図 9 に示す例では、集合 A におけるハッシュ関数 h_1 の最小ハッシュ値は 1 2 であり、ハッシュ関数 h_2 の最小ハッシュ値は 5 6 である。このように、最小ハッシュ値特定部 3 は、集合 A に関して、ハッシュ関数毎に最小ハッシュ値を特定する。さらに、最小ハッシュ値特定部 3 は、他の各集合に関しても、同様に、ハッシュ関数毎に最小ハッシュ値を特定する。

【 0 0 4 5 】

図 4 から図 9 では、複数の集合が 2 つの集合 A、B である場合を例にして説明した。類似度導出システム 1 に与えられる集合が 3 つ以上の場合であっても、第 1 のハッシュ値計算部 5 1、全体集合生成部 5 2、インデックス情報生成部 5 3、第 2 のハッシュ値計算部 5 4、ハッシュ値決定部 5 5、および、最小ハッシュ値特定部 3 の動作は、上記の動作と同様である。

20

【 0 0 4 6 】

第 1 のハッシュ値計算部 5 1、全体集合生成部 5 2、インデックス情報生成部 5 3、第 2 のハッシュ値計算部 5 4、ハッシュ値決定部 5 5 を含むハッシュ値計算部 2 は、例えば、類似度導出プログラムに従って動作するコンピュータの CPU によって実現される。例えば、CPU が、コンピュータのプログラム記憶装置等のプログラム記録媒体から類似度導出プログラムを読み込み、その類似度導出プログラムに従って、第 1 のハッシュ値計算部 5 1、全体集合生成部 5 2、インデックス情報生成部 5 3、第 2 のハッシュ値計算部 5 4、ハッシュ値決定部 5 5 を含むハッシュ値計算部 2 として動作すればよい。

30

【 0 0 4 7 】

次に、第 2 の実施形態の処理経過の例を示す。なお、既に説明した事項については、適宜、説明を省略する。図 10 は、第 2 の実施形態の処理経過の例を示すフローチャートである。また、上記の例と同様に、複数のハッシュ関数を h_1, h_2, \dots, h_m とし、 h_1 が所定のハッシュ関数であるものとする。

【 0 0 4 8 】

まず、第 1 のハッシュ値計算部 5 1 が、複数の集合に含まれる各集合の個々の要素に対して所定のハッシュ関数 h_1 を適用することによって、第 1 のハッシュ値を計算する（ステップ S 5 1）。この結果、各集合の個々の要素毎に、第 1 のハッシュ値が得られる。

40

【 0 0 4 9 】

次に、全体集合生成部 5 2 が、全体集合を生成する（ステップ S 5 2）。ステップ S 5 2 において、全体集合生成部 5 2 は、各集合の個々の要素のうち、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、その複数の要素に該当しない各要素を取り出し、取り出した各要素を含む 1 つの集合を全体集合とすればよい。

【 0 0 5 0 】

次に、インデックス情報生成部 5 3 が、インデックス情報を生成する（ステップ S 5 3

50

）。例えば、インデックス情報生成部 5 3 は、全体集合に属する各要素に対して識別情報を割り当てる。そして、インデックス情報生成部 5 3 は、複数の集合に含まれる各集合の各要素毎に、その要素に該当する全体集合内の要素の識別情報を定めた情報を、インデックス情報として生成する。

【 0 0 5 1 】

次に、第 2 のハッシュ値計算部 5 4 が、全体集合に属する各要素に対して、所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値をそれぞれ計算する（ステップ S 5 4）。

【 0 0 5 2 】

次に、ハッシュ値決定部 5 5 が、インデックス情報に基づいて、各集合の個々の要素の第 1 のハッシュ値と、ステップ S 5 4 で計算された各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値とを組み合わせることによって、各集合の個々の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値を決定する（ステップ S 5 5）。ステップ S 5 の結果、例えば、図 8 に例示するように、各集合の個々の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する複数のハッシュ値が得られる。

10

【 0 0 5 3 】

ステップ S 5 1 ~ ステップ S 5 5 は、第 1 の実施形態のステップ S 1 をより具体化した処理の一例である。

【 0 0 5 4 】

ステップ S 5 5 の次に、最小ハッシュ値特定部 3 が、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に最小ハッシュ値を特定する（ステップ S 2）。ステップ S 2 の結果、例えば、図 9 に例示するように、各集合に関して、ハッシュ関数毎に最小ハッシュ値が特定される。

20

【 0 0 5 5 】

次に、類似度導出部 4 が、複数の集合から得られる 1 組以上の集合のペアに関して、ペアをなす 2 つの集合の類似度を導出する（ステップ S 3）。

【 0 0 5 6 】

ステップ S 2, S 3 は、第 1 の実施形態におけるステップ S 2, S 3 と同様である。

【 0 0 5 7 】

以下、ステップ S 3 についてより具体的に説明する。図 1 1 は、ステップ S 3 の動作の具体例の一例を示すフローチャートである。なお、類似度導出システム 1 に与えられる集合の数は、2 つであっても、3 つ以上であってもよい。以下の説明では、類似度導出システム 1 が、与えられた集合から得られる全てのペア（集合のペア）に関して、ペアをなす 2 つの集合の類似度を算出する場合を示す。ただし、類似度導出システム 1 に与えられた集合から得られる全てのペアのうち、例えば、類似度導出システム 1 の操作者によって指定されたペアに関してのみ、類似度を算出してもよい。

30

【 0 0 5 8 】

類似度導出部 4 は、複数の集合から、集合のペアを 1 つ取り出す（ステップ S 6 1）。1 つのペアは、2 つの集合からなる。また、ステップ S 6 1 では、類似度導出部 4 は、まだ選択されていないペアを 1 つ取り出す。

40

【 0 0 5 9 】

類似度導出部 4 は、ステップ S 6 1 で選択されたペアをなす 2 つの集合の類似度を、その 2 つの集合における個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する（ステップ S 6 2）。例えば、類似度導出部 4 は、ペアをなす 2 つの集合に関して、個々のハッシュ値に対応する最小ハッシュ値同士が一致している数（以下、一致数と記す。）を、ペアをなす 2 つの集合の類似度として定めてもよい。また、例えば、類似度導出部 4 は、ハッシュ関数の数に対する一致数の割合を、類似度として定めてもよい。

【 0 0 6 0 】

ステップ S 6 2 の次に、類似度導出部 4 は、未選択の集合のペアが存在するか否かを判

50

定する（ステップ S 6 3）。

【 0 0 6 1 】

未選択の集合のペアが存在するならば（ステップ S 6 3 の Yes）、類似度導出部 4 は、ステップ S 6 1 以降の動作を繰り返す。

【 0 0 6 2 】

未選択の集合のペアが存在しないならば（ステップ S 6 3 の No）、全てのペアに関して、ペアをなす 2 つの集合の類似度が導出されていることになるので、その時点で処理を終了する。

【 0 0 6 3 】

本実施形態によれば、全体集合生成部 5 2 が、各集合の個々の要素のうち、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、その複数の要素に該当しない各要素を取り出し、取り出した各要素を含む 1 つの全体集合を生成する。そして、第 2 のハッシュ値計算部 5 4 が、全体集合に属する各要素に対して、所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値をそれぞれ計算する。従って、第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素については、第 2 のハッシュ値計算部 5 4 は、ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の計算を 1 回のみ行う。そのような複数の要素それぞれに対して、ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の計算を行っても、同一の計算を行うことになり、同一のハッシュ値が得られるが、上記のように、第 2 のハッシュ値計算部 5 4 は、そのような複数の要素については、ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値の計算を 1 回のみ行う。従って、同一の計算を重複して行うことがなく、集合の類似度を導出する際に、ハッシュ値の計算量を低減させることができる。

【 0 0 6 4 】

実施形態 3 .

第 3 の実施形態は、第 1 の実施形態をより具体的に示した実施形態である。図 1 2 は、本発明の第 3 の実施形態の類似度導出システムの例を示すブロック図である。本実施形態の類似度導出システム 1 も、ハッシュ値計算部 2 と、最小ハッシュ値特定部 3 と、類似度導出部 4 とを備える。最小ハッシュ値特定部 3 および類似度導出部 4 は、第 1 の実施形態および第 2 の実施形態における最小ハッシュ値特定部 3 および類似度導出部 4 と同様である。

【 0 0 6 5 】

また、本実施の形態においても、第 2 の実施形態と同様に、複数のハッシュ関数の数は、 m 個であり、個々のハッシュ関数を h_1, h_2, \dots, h_m と記すこととする。そして、 h_1 が、所定のハッシュ関数であるものとする。

【 0 0 6 6 】

既に説明したように、複数のハッシュ関数の中で値域が最も広いハッシュ関数を所定のハッシュ関数として定めることが好ましい。すなわち、所定のハッシュ関数は、複数のハッシュ関数の中で最も値域が広いハッシュ関数であることが好ましい。本例では、ハッシュ関数 h_1 の値域が、複数のハッシュ関数 $h_1 \sim h_m$ の値域の中で最も広いものとする。

【 0 0 6 7 】

また、所定のハッシュ関数 h_1 で計算されたハッシュ値を、第 1 のハッシュ値と記す。

【 0 0 6 8 】

第 3 の実施形態のハッシュ値計算部 2 は、集合選択部 6 1 と、要素選択部 6 2 と、第 1 のハッシュ値計算部 6 3 と、判定部 6 4 と、第 2 のハッシュ値計算部 6 5 とを備える。なお、第 3 の実施形態における第 1 のハッシュ値計算部 6 3 および第 2 のハッシュ値計算部 6 5 の動作は、第 2 の実施形態における第 1 のハッシュ値計算部 5 1 および第 2 のハッシュ値計算部 5 4 の動作とは異なる。

【 0 0 6 9 】

集合選択部 6 1 は、複数の集合から、順次、1 つの集合を選択する。

【 0 0 7 0 】

10

20

30

40

50

要素選択部 6 2 は、集合選択部 6 1 によって選択された集合から、順次、1 つの要素を選択する。

【 0 0 7 1 】

第 1 のハッシュ値計算部 6 3 は、要素選択部 6 2 によって選択された要素に所定のハッシュ関数を適用することによって、その選択された要素の第 1 のハッシュ値を計算する。

【 0 0 7 2 】

判定部 6 4 は、要素選択部 6 2 によって選択された要素と、第 1 のハッシュ値が一致し、かつ、要素自体が一致する要素が既に選択されているか否かを判定する。

【 0 0 7 3 】

以下、要素選択部 6 2 によって選択された要素と、第 1 のハッシュ値が一致し、かつ、要素自体が一致する要素を、一致要素と記す。

10

【 0 0 7 4 】

第 2 のハッシュ値計算部 6 5 は、一致要素が既に選択されている場合に、要素選択部 6 2 によって選択された要素の所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を、その一致要素の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値と同一であると定める。

【 0 0 7 5 】

また、第 2 のハッシュ値計算部 6 5 は、一致要素が選択されていない場合に、要素選択部 6 2 によって選択された要素の所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を計算する。

20

【 0 0 7 6 】

集合選択部 6 1、要素選択部 6 2、第 1 のハッシュ値計算部 6 3、判定部 6 4、第 2 のハッシュ値計算部 6 5 を含むハッシュ値計算部 2 は、例えば、類似度導出プログラムに従って動作するコンピュータの CPU によって実現される。例えば、CPU が、コンピュータのプログラム記憶装置等のプログラム記録媒体から類似度導出プログラムを読み込み、その類似度導出プログラムに従って、集合選択部 6 1、要素選択部 6 2、第 1 のハッシュ値計算部 6 3、判定部 6 4、第 2 のハッシュ値計算部 6 5 を含むハッシュ値計算部 2 として動作すればよい。

【 0 0 7 7 】

次に、第 3 の実施形態の処理経過の例を示す。なお、既に説明した事項については、適宜、説明を省略する。図 1 3 および図 1 4 は、第 3 の実施形態の処理経過の例を示すフローチャートである。

30

【 0 0 7 8 】

まず、集合選択部 6 1 が、複数の集合から 1 つの集合を選択する (ステップ S 7 1)。ステップ S 7 1 で、集合選択部 6 1 は、まだ選択されていない集合を 1 つ選択する。

【 0 0 7 9 】

本例では、図 1 5 に示す集合が選択されたものとする。以下、図 1 5 に示す集合を集合 A と記す。なお、図 1 5 は、複数の集合のうちの 1 つの集合の例を示す模式図である。

【 0 0 8 0 】

次に、要素選択部 6 2 が、ステップ S 7 1 で選択された集合から 1 つの要素を選択する (ステップ S 7 2)。ステップ S 7 2 で、要素選択部 6 2 は、まだ選択されていない要素を 1 つ選択する。ここでは、集合 A から要素 "the" を選択したものとする。

40

【 0 0 8 1 】

次に、第 1 のハッシュ値計算部 6 3 が、ステップ S 7 2 で選択された要素に所定のハッシュ関数 h_1 を適用することによって、その要素の第 1 のハッシュ値を計算する (ステップ S 7 3)。図 1 6 は、最初にステップ S 7 3 を実行した後に得られたハッシュ値の例を示す模式図である。図 1 6 に示すように、この時点では、要素 "the" の第 1 のハッシュ値が得られているが、他のハッシュ値は得られていない。

【 0 0 8 2 】

ステップ S 7 3 の後に、判定部 6 4 は、ステップ S 7 2 で選択された要素と、第 1 のハ

50

ッシュ値が一致し、かつ、要素自体が一致する要素（すなわち、一致要素）が既に選択されているか否かを判定する（ステップS 7 4）。

【0083】

一致要素が既に選択されている場合（ステップS 7 4のYes）、ステップS 7 5に移行し、一致要素が選択されていない場合（ステップS 7 4のNo）、ステップS 7 6に移行する。本例では、集合Aの要素"the"は、1番目に選択された要素であり、一致要素は選択されていない。したがって、ステップS 7 6に移行する。

【0084】

ステップS 7 6では、第2のハッシュ値計算部65が、ステップS 7 2で選択された要素の所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を計算する。図17は、最初にステップS 7 6を実行した後に得られたハッシュ値の例を示す模式図である。図17に示すように、この時点では、要素"the"の複数のハッシュ関数 h_1, h_2, \dots, h_m に対応するハッシュ値が得られている。

10

【0085】

ステップS 7 5またはステップS 7 6の実行後には、要素選択部62が、選択されている集合の中に未選択の要素があるか否かを判定する（ステップS 7 7）。本例では、選択されている集合Aの中に未選択の要素があると判定する（ステップS 7 7のYes）。この場合、ステップS 7 2以降の処理を繰り返す。

【0086】

集合Aの要素は全て異なる要素であるので（図15参照）、集合Aを選択しているときにステップS 7 4では、一致要素は選択されていないと判定され、ステップS 7 6に移行する。図18は、ステップS 7 7で未選択の要素がないと判定されるまでに得られたハッシュ値の例を示す模式図である。

20

【0087】

ステップS 7 7で、選択している集合の中に未選択の要素がないと判定した場合（ステップS 7 7のNo）、ステップS 7 8に移行する。ステップS 7 8では、集合選択部61が、未選択の集合があるか否かを判定する。

【0088】

未選択の集合があると判定した場合（ステップS 7 8のYes）、ステップS 7 1以降の処理を繰り返す。ここでは、まだ、図19に示す集合Bが選択されていないものとする。なお、図19は、複数の集合のうちの1つの集合の例を示す模式図である。

30

【0089】

ステップS 7 1で集合選択部61が集合Bを選択し、ステップS 7 2で要素選択部62が集合Bの要素"room"を選択したとする。その後、ステップS 7 3が実行される。この場合、要素"room"に応じた一致要素は選択されていないので（ステップS 7 4のNo）、ステップS 7 6に移行する。この結果、集合Bの要素"room"の複数のハッシュ関数 h_1, h_2, \dots, h_m に対応するハッシュ値が得られる。

【0090】

次に、ステップS 7 2で、要素選択部62が、集合Bの要素"mountain"を選択したとする。すると、第1のハッシュ値計算部63は、集合Bの要素"mountain"に所定のハッシュ関数を適用することによって、その要素の第1のハッシュ値を計算する（ステップS 7 3）。図20は、集合Bの要素"mountain"が選択され、ステップS 7 3を実行した後に得られたハッシュ値の例を示す模式図である。

40

【0091】

次に、判定部64は、ステップS 7 2で選択された要素（ここでは、集合Bの要素"mountain"）と、第1のハッシュ値が一致し、かつ、要素自体が一致する要素（すなわち、一致要素）が既に選択されているか否かを判定する（ステップS 7 4）。

【0092】

この時点では、集合Bの要素"mountain"と、第1のハッシュ値が一致し、かつ、要素自体が一致する要素である「集合Aの要素"mountain"」が選択済みである（ステップS

50

74のYes)。従って、ステップS75に移行する。なお、この場合、「集合Aの要素"mountain"」が一致要素に該当する。

【0093】

ステップS75では、第2のハッシュ値計算部65は、ステップS72で選択された「集合Bの要素"mountain"」の所定のハッシュ関数 h_1 以外の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を、一致要素(集合Aの要素"mountain")の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値と同一であると定める。図21は、集合Bの要素"mountain"を選択したときにおけるステップS75の実行後に得られているハッシュ値の例を示す模式図である。

【0094】

未選択の要素"view"が存在するので(ステップS77のYes)、ステップS72で、要素選択部62が、集合Bの要素"view"を選択し、ステップS73以降の処理を繰り返す。この場合、要素"view"に応じた一致要素は選択されていないので(ステップS74のNo)、ステップS76に移行する。この結果、集合Bの要素"view"の複数のハッシュ関数 h_1, h_2, \dots, h_m に対応するハッシュ値が得られる。

【0095】

ここで、集合Bの要素は全て選択された状態になっているので(ステップS77のNo)、ステップS78に移行する。図22は、この時点で得られているハッシュ値の例を示す模式図である。

【0096】

ステップS78で、判定部64が、未選択の集合がないと判定したとする(ステップS78のNo)。この場合、ステップS2(図14参照)に移行する。

【0097】

ステップS2に移行するまでの処理(ステップS78で未選択の集合がないと判定するまでの処理)は、第1の実施形態のステップS1をより具体化した処理の一例である。

【0098】

ステップS2では、最小ハッシュ値特定部3が、複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に最小ハッシュ値を特定する。本例では、複数の集合が、2つの集合A, Bであるとする。この場合、ステップS2の結果、第2の実施形態で示した図9に例示するように、各集合に関して、ハッシュ関数毎に最小ハッシュ値が特定される。

【0099】

ステップS2の次に、類似度導出部4が、複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を導出する(ステップS3)。

【0100】

ステップS2, S3は、第1の実施形態および第2の実施形態におけるステップS2, S3と同様である。

【0101】

ステップS3のより具体的な動作の例については、第2の実施形態で図11を参照して説明したので、ここでは説明を省略する。

【0102】

本実施形態によれば、判定部64が、ステップS72で選択された要素に応じた一致要素が既に選択済みであると判定した場合(ステップS74のYes)、ステップS75に移行する。一致要素の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値は、既に計算されている。そして、ステップS75において、第2のハッシュ値計算部65は、ステップS72で選択された要素の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を、一致要素の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値と同一であると定める。ステップS72で選択された要素に応じた一致要素は、ステップS72で選択された要素と一致する。従って、ステップS75において、第2のハッシュ値計算部65は、ステップS72で選択された要素の各ハッシュ関数 h_2, \dots, h_m に対応するハ

10

20

30

40

50

ハッシュ値を、一致要素の各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値と同一であると定めることができる。そして、ステップ S 7 2 で選択された要素に応じた一致要素が既に選択済みである場合に、ステップ S 7 5 に移行する。その結果、既に計算された各ハッシュ関数 h_2, \dots, h_m に対応するハッシュ値を再度計算する必要がない。よって、集合の類似度を導出する際に、ハッシュ値の計算量を低減させることができる。

【0103】

次に、各実施形態の変形例について説明する。この変形例では、複数の集合に含まれる各集合の各要素を、複数のハッシュ関数 h_1, h_2, \dots, h_m とは別のハッシュ関数 (h_0 と記す。) によって、ハッシュ値 (数値) に変換する。このハッシュ関数 h_0 を数値化ハッシュ関数と記す。例えば、各集合の各要素が文字列である場合、それぞれの文字列に数値化ハッシュ関数 h_0 を適用して、各集合の各要素を文字列から数値に変換する。

10

【0104】

そして、変換後の各要素を、各集合の各要素として、上記の第 1 の実施形態、第 2 の実施形態、第 3 の実施形態を適用してもよい。文字列に対して複数のハッシュ関数のハッシュ値を計算する計算量よりも、数値に対して複数のハッシュ関数のハッシュ値を計算する計算量の方が少ない。従って、上記のように、各集合の各要素を文字列から数値に変換することで計算量を少なくすることができる。

【0105】

ただし、数値化ハッシュ関数 h_0 による変換後の各要素は、数値である。この場合、第 1 の実施形態のステップ S 1 では、ハッシュ値計算部 2 は、数値化ハッシュ関数 h_0 による変換後のハッシュ値が一致する複数の要素に関しては、複数のハッシュ関数の計算の重複を排除し、各集合の個々の要素に対して複数のハッシュ値を求めればよい。

20

【0106】

また、本変形例を第 2 の実施形態に適用する場合には、第 1 のハッシュ値計算部 5 1 が、数値化ハッシュ関数 h_0 を用いて、各集合の各要素を数値に変換する。そして、全体集合生成部 5 2 は、各集合の個々の要素 (変換後の要素) のうち、数値化ハッシュ関数 h_0 による変換後のハッシュ値が一致する複数の要素からは 1 つの要素だけを取り出すとともに、その複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を 1 つ含む 1 つの全体集合を生成すればよい。また、この場合、第 2 のハッシュ値計算部 5 4 は、全体集合に属する各要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応するハッシュ値をそれぞれ計算する。ハッシュ値決定部 5 5 は、インデックス情報に基づいて、各集合の個々の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応するハッシュ値を特定する。

30

【0107】

また、本変形例を第 3 の実施形態に適用する場合には、ステップ S 7 3 において、第 1 のハッシュ値計算部 6 3 が、ステップ S 7 2 で選択された要素に数値化ハッシュ関数 h_0 を適用して、その要素を数値の要素に変換する。ステップ S 7 4 では、判定部 6 4 は、数値化ハッシュ関数 h_0 で変換された要素と一致する要素を一致要素とし、一致要素が既に得られているか否かを判定すればよい。そして、一致要素が得られていない場合、ステップ S 7 6 で、第 2 のハッシュ値計算部 6 5 は、変換後の要素に対して、複数のハッシュ関数 h_1, h_2, \dots, h_m を適用することによって、複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する各ハッシュ値を計算する。また、一致要素が得られている場合、ステップ S 7 5 で、第 2 のハッシュ値計算部 6 5 は、変換後の要素の複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する各ハッシュ値を、一致要素の複数のハッシュ関数 h_1, h_2, \dots, h_m に対応する各ハッシュ値と同一であると定める。

40

【0108】

図 2 3 は、本発明の実施形態の類似度導出システム 1 に係るコンピュータの構成例を示す概略ブロック図である。コンピュータ 1 0 0 0 は、CPU 1 0 0 1 と、主記憶装置 1 0 0 2 と、補助記憶装置 1 0 0 3 と、インタフェース 1 0 0 4 とを備える。

【0109】

50

本発明の各実施形態の類似度導出システム 1 は、例えば、コンピュータ 1000 によって実現される。類似度導出システム 1 の動作は、類似度導出プログラムの形式で補助記憶装置 1003 に記憶されている。CPU 1001 は、その類似度導出プログラムを読み出し、類似度導出プログラムを主記憶装置 1002 に展開し、その類似度導出プログラムに従って、上記の各実施形態で説明した処理を実行する。

【0110】

補助記憶装置 1003 は、一時的でない有形の媒体の例である。一時的でない有形の媒体の他の例として、インタフェース 1004 を介して接続される磁気ディスク、光磁気ディスク、CD-ROM (Compact Disk Read Only Memory)、DVD-ROM (Digital Versatile Disk Read Only Memory)、半導体メモリ等が挙げられる。また、プログラムが通信回線によってコンピュータ 1000 に配信される場合、配信を受けたコンピュータ 1000 がそのプログラムを主記憶装置 1002 に展開し、そのプログラムに従って上記の実施形態で説明した処理を実行してもよい。

10

【0111】

また、各構成要素の一部または全部は、汎用または専用の回路 (circuitry)、プロセッサ等やこれらの組合せによって実現されてもよい。これらは、単一のチップによって構成されてもよいし、バスを介して接続される複数のチップによって構成されてもよい。各構成要素の一部または全部は、上述した回路等とプログラムとの組合せによって実現されてもよい。

【0112】

各構成要素の一部または全部が複数の情報処理装置や回路等により実現される場合には、複数の情報処理装置や回路等は集中配置されてもよいし、分散配置されてもよい。例えば、情報処理装置や回路等は、クライアントアンドサーバシステム、クラウドコンピューティングシステム等、各々が通信ネットワークを介して接続される形態として実現されてもよい。

20

【0113】

上記の本発明の実施形態は、以下の付記のようにも記載され得るが、以下に限定されるわけではない。

【0114】

(付記 1)

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算手段と、

30

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定手段と、

前記複数の集合から得られる 1 組以上の集合のペアに関して、ペアをなす 2 つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出手段とを備える

40

ことを特徴とする類似度導出システム。

【0115】

(付記 2)

前記ハッシュ値計算手段は、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第 1 のハッシュ値を計算する第 1 のハッシュ値計算手段と、

前記各集合の個々の要素のうち、前記第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む 1 つの集合である全体集合を生成する全体集合生成手段と、

50

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成手段と、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算手段と、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第1のハッシュ値と、当該第1のハッシュ値に対応する、前記第2のハッシュ値計算手段によって計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定手段とを含む

付記1に記載の類似度導出システム。

10

【0116】

(付記3)

前記ハッシュ値計算手段は、

前記複数の集合から、順次、1つの集合を選択する集合選択手段と、

選択された集合から、順次、1つの要素を選択する要素選択手段と、

選択された要素に前記所定のハッシュ関数を適用することによって前記選択された要素の第1のハッシュ値を計算する第1のハッシュ値計算手段と、

前記選択された要素と、第1のハッシュ値が一致し、かつ、要素自体が一致する要素である一致要素が既に選択されているか否かを判定する判定手段と、

前記一致要素が既に選択されている場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を、前記一致要素の前記各ハッシュ関数に対応するハッシュ値と同一であると定め、

20

前記一致要素が選択されていない場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を計算する第2のハッシュ値計算手段とを含む

付記1に記載の類似度導出システム。

【0117】

(付記4)

前記所定のハッシュ関数は、前記複数のハッシュ関数の中で最も値域が広いハッシュ関数である

30

付記1から付記3のうちのいずれか1項に記載の類似度導出システム。

【0118】

(付記5)

前記類似度導出手段は、

個々のハッシュ関数に対応する最小ハッシュ値同士が一致している数を、ペアをなす2つの集合の類似度として定める

付記1から付記4のうちのいずれか1項に記載の類似度導出システム。

【0119】

(付記6)

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算処理と、

40

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理と、

前記複数の集合から得られる1組以上の集合のペアに関して、ペアをなす2つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理とを含む

ことを特徴とする類似度導出方法。

50

【 0 1 2 0 】

(付記 7)

前記ハッシュ値計算処理は、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第 1 のハッシュ値を計算する第 1 のハッシュ値計算処理と、

前記各集合の個々の要素のうち、前記第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む 1 つの集合である全体集合を生成する全体集合生成処理と、

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成処理と、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第 2 のハッシュ値計算処理と、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第 1 のハッシュ値と、当該第 1 のハッシュ値に対応する、前記第 2 のハッシュ値計算処理で計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定処理とを含む

付記 6 に記載の類似度導出方法。

【 0 1 2 1 】

(付記 8)

前記ハッシュ値計算処理は、

前記複数の集合から、順次、1 つの集合を選択する集合選択処理と、

選択された集合から、順次、1 つの要素を選択する要素選択処理と、

選択された要素に前記所定のハッシュ関数を適用することによって前記選択された要素の第 1 のハッシュ値を計算する第 1 のハッシュ値計算処理と、

前記選択された要素と、第 1 のハッシュ値が一致し、かつ、要素自体が一致する要素である一致要素が既に選択されているか否かを判定する判定処理と、

前記一致要素が既に選択されている場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を、前記一致要素の前記各ハッシュ関数に対応するハッシュ値と同一であると定め、

前記一致要素が選択されていない場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を計算する第 2 のハッシュ値計算処理とを含む

付記 6 に記載の類似度導出方法。

【 0 1 2 2 】

(付記 9)

コンピュータに、

複数の集合に含まれる各集合の個々の要素に対して複数のハッシュ関数を適用して得られる複数のハッシュ値を求める際に、前記複数のハッシュ関数のうちの所定のハッシュ関数によって得られるハッシュ値が一致し、かつ、要素自体が一致する複数の要素に関しては、前記所定のハッシュ関数以外の各ハッシュ関数の計算の重複を排除し、前記各集合の個々の要素に対して前記複数のハッシュ値を求めるハッシュ値計算処理、

前記複数の集合に含まれるそれぞれの集合に関して、個々のハッシュ関数毎に、ハッシュ値の最小値である最小ハッシュ値を特定する最小ハッシュ値特定処理、および、

前記複数の集合から得られる 1 組以上の集合のペアに関して、ペアをなす 2 つの集合の類似度を、個々のハッシュ関数に対応する最小ハッシュ値に基づいて導出する類似度導出処理

を実行させるための類似度導出プログラムを記録したコンピュータ読取可能な記録媒体。

【 0 1 2 3 】

10

20

30

40

50

(付記 1 0)

前記コンピュータに、

前記ハッシュ値計算処理で、

前記各集合の個々の要素に対して前記所定のハッシュ関数を適用することによって第 1 のハッシュ値を計算する第 1 のハッシュ値計算処理、

前記各集合の個々の要素のうち、前記第 1 のハッシュ値が一致し、かつ、要素自体が一致する複数の要素からは 1 つの要素だけを取り出すとともに、前記複数の要素に該当しない各要素を取り出すことによって、重複なく全ての種類の要素を含む 1 つの集合である全体集合を生成する全体集合生成処理、

前記複数の集合に含まれるどの集合のどの要素が前記全体集合のどの要素に該当するかを示すインデックス情報を生成するインデックス情報生成処理、

前記全体集合に属する各要素に対して、前記複数のハッシュ関数のうちの前記所定のハッシュ関数以外の各ハッシュ関数を適用することによって、前記各ハッシュ関数に対応するハッシュ値を計算する第 2 のハッシュ値計算処理、および、

前記インデックス情報に基づいて、前記各集合の個々の要素の前記第 1 のハッシュ値と、当該第 1 のハッシュ値に対応する、前記第 2 のハッシュ値計算処理で計算されたハッシュ値とを組み合わせることによって、前記各集合の個々の要素に対して、前記複数のハッシュ関数に対応する複数のハッシュ値を決定するハッシュ値決定処理を実行させる

類似度導出プログラムを記録した付記 9 に記載のコンピュータ読取可能な記録媒体。

【 0 1 2 4 】

(付記 1 1)

前記コンピュータに、

前記ハッシュ値計算処理で、

前記複数の集合から、順次、1 つの集合を選択する集合選択処理、

選択された集合から、順次、1 つの要素を選択する要素選択処理、

選択された要素に前記所定のハッシュ関数を適用することによって前記選択された要素の第 1 のハッシュ値を計算する第 1 のハッシュ値計算処理、

前記選択された要素と、第 1 のハッシュ値が一致し、かつ、要素自体が一致する要素である一致要素が既に選択されているか否かを判定する判定処理、および、

前記一致要素が既に選択されている場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を、前記一致要素の前記各ハッシュ関数に対応するハッシュ値と同一であると定め、

前記一致要素が選択されていない場合に、前記選択された要素の前記所定のハッシュ関数以外の各ハッシュ関数に対応するハッシュ値を計算する第 2 のハッシュ値計算処理を実行させる

類似度導出プログラムを記録した付記 9 に記載のコンピュータ読取可能な記録媒体。

【 0 1 2 5 】

以上、実施形態を参照して本願発明を説明したが、本願発明は上記の実施形態に限定されるものではない。本願発明の構成や詳細には、本願発明の範囲内で当業者が理解し得る様々な変更をすることができる。

【 産業上の利用の可能性 】

【 0 1 2 6 】

本発明は、集合同士の類似度を導出する類似度導出システムに好適に適用される。

【 符号の説明 】

【 0 1 2 7 】

1 類似度導出システム

2 ハッシュ値計算部

3 最小ハッシュ値特定部

4 類似度導出部

5 1 第 1 のハッシュ値計算部

10

20

30

40

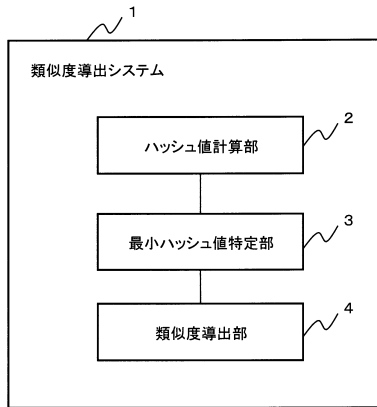
50

- 5 2 全体集合生成部
- 5 3 インデックス情報生成部
- 5 4 第2のハッシュ値計算部
- 5 5 ハッシュ値決定部
- 6 1 集合選択部
- 6 2 要素選択部
- 6 3 第1のハッシュ値計算部
- 6 4 判定部
- 6 5 第2のハッシュ値計算部

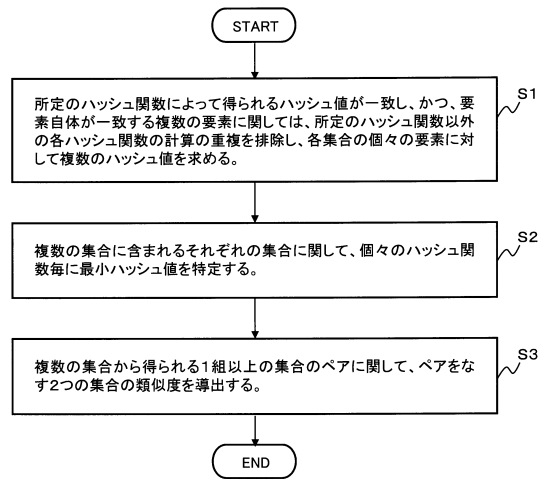
【図面】

10

【図 1】



【図 2】



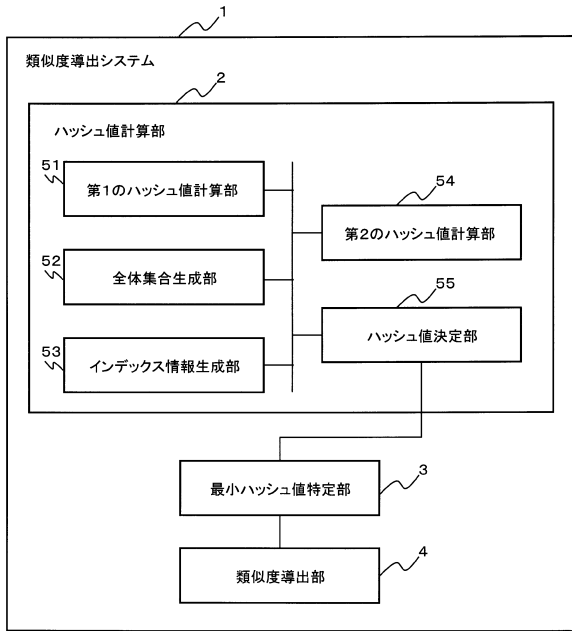
20

30

40

50

【図3】



【図4】

集合A

要素	h1
"the"	389
"highest"	12
"mountain"	666
"Fuji"	920

集合B

要素	h1
"room"	124
"mountain"	666
"view"	300

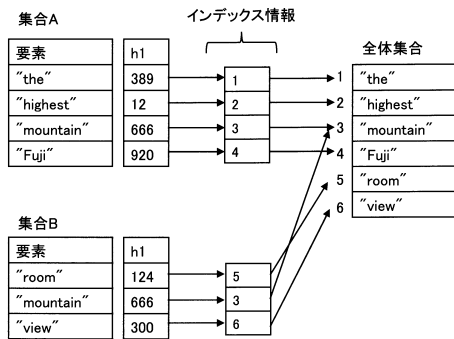
10

【図5】

全体集合

"the"
"highest"
"mountain"
"Fuji"
"room"
"view"

【図6】



20

30

40

50

【 図 7 】

全体集合	h2	h3	...	hm
"the"	56	143	...	513
"highest"	323	902	...	322
"mountain"	92	792	...	30
"Fuji"	820	43	...	325
"room"	329	469	...	832
"view"	810	521	...	292

【 図 8 】

集合A

要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"	12	323	902	...	322
"mountain"	666	92	792	...	30
"Fuji"	920	820	43	...	325

集合B

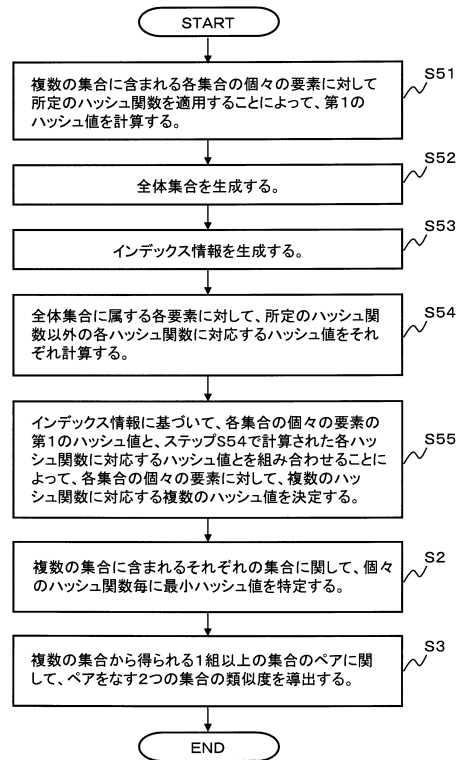
要素	h1	h2	h3	...	hm
"room"	124	329	469	...	832
"mountain"	666	92	792	...	30
"view"	300	810	521	...	292

10

【 図 9 】

	h1	h2	h3	...	hm
集合Aの最小ハッシュ値	12	56	43	...	30
集合Bの最小ハッシュ値	124	92	469	...	30

【 図 10 】



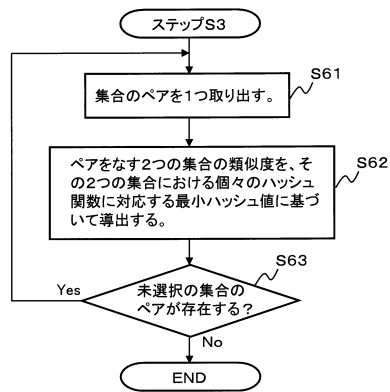
20

30

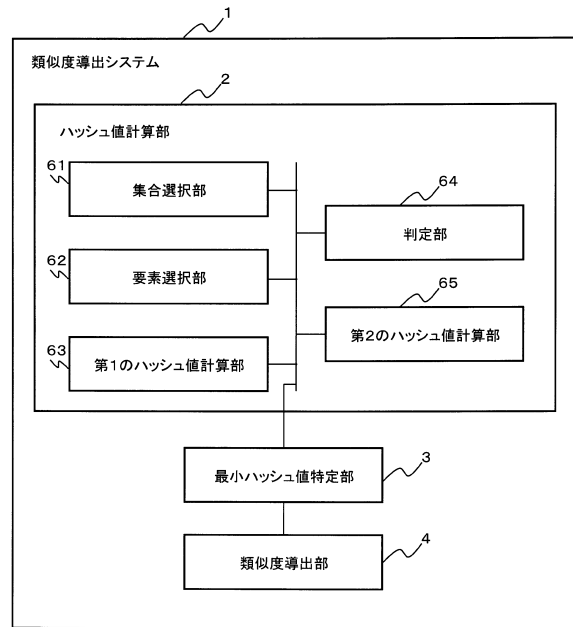
40

50

【図 1 1】



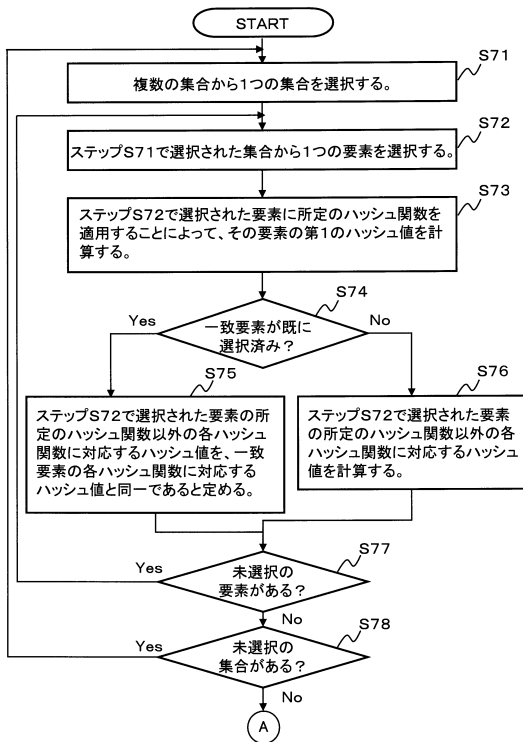
【図 1 2】



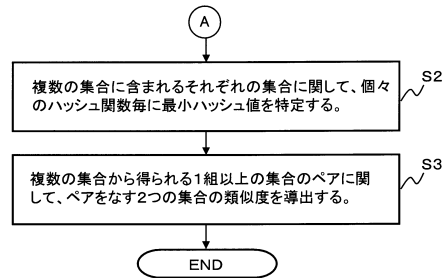
10

20

【図 1 3】



【図 1 4】



30

40

50

【図 1 5】

集合A

要素
"the"
"highest"
"mountain"
"Fuji"

【図 1 6】

集合A

要素	h1	h2	h3	...	hm
"the"	389				
"highest"					
"mountain"					
"Fuji"					

10

【図 1 7】

集合A

要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"					
"mountain"					
"Fuji"					

【図 1 8】

集合A

要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"	12	323	902	...	322
"mountain"	666	92	792	...	30
"Fuji"	920	820	43	...	325

20

【図 1 9】

集合B

要素
"room"
"mountain"
"view"

【図 2 0】

集合A

要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"	12	323	902	...	322
"mountain"	666	92	792	...	30
"Fuji"	920	820	43	...	325

集合B

要素	h1	h2	h3	...	hm
"room"	124	329	469	...	832
"mountain"	666				
"view"					

30

40

50

【図 2 1】

集合A

要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"	12	323	902	...	322
"mountain"	666	92	792	...	30
"Fuji"	920	820	43	...	325

集合B

要素	h1	h2	h3	...	hm
"room"	124	329	469	...	832
"mountain"	666	92	792	...	30
"view"					

【図 2 2】

集合A

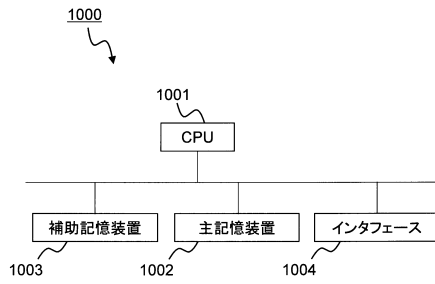
要素	h1	h2	h3	...	hm
"the"	389	56	143	...	513
"highest"	12	323	902	...	322
"mountain"	666	92	792	...	30
"Fuji"	920	820	43	...	325

集合B

要素	h1	h2	h3	...	hm
"room"	124	329	469	...	832
"mountain"	666	92	792	...	30
"view"	300	810	521	...	292

10

【図 2 3】



20

30

40

50

フロントページの続き

- (56)参考文献 国際公開第2021/038887(WO, A1)
米国特許出願公開第2018/0181609(US, A1)
米国特許出願公開第2018/0095941(US, A1)
米国特許出願公開第2017/0322930(US, A1)
米国特許出願公開第2017/0161375(US, A1)
米国特許出願公開第2017/0078286(US, A1)
- (58)調査した分野 (Int.Cl., DB名)
G06F 16/00 - 16/958
G09C 1/00