

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5258140号
(P5258140)

(45) 発行日 平成25年8月7日(2013.8.7)

(24) 登録日 平成25年5月2日(2013.5.2)

(51) Int. Cl.	F I
G06F 17/30 (2006.01)	G06F 17/30 220Z
	G06F 17/30 340A
	G06F 17/30 340B
	G06F 17/30 350C

請求項の数 11 (全 17 頁)

(21) 出願番号	特願2003-544641 (P2003-544641)	(73) 特許権者	590000248
(86) (22) 出願日	平成14年10月22日 (2002.10.22)		コーニンクレッカ フィリップス エレク
(65) 公表番号	特表2005-509954 (P2005-509954A)		トロニクス エヌ ヴィ
(43) 公表日	平成17年4月14日 (2005.4.14)		オランダ国 5656 アーエー アイ
(86) 国際出願番号	PCT/IB2002/004413		ドーフエン ハイテック キャンパス 5
(87) 国際公開番号	W02003/042879	(74) 代理人	100087789
(87) 国際公開日	平成15年5月22日 (2003.5.22)		弁理士 津軽 進
審査請求日	平成17年10月21日 (2005.10.21)	(74) 代理人	100122769
審査番号	不服2010-16461 (P2010-16461/J1)		弁理士 笛田 秀仙
審査請求日	平成22年7月22日 (2010.7.22)	(74) 代理人	100163821
(31) 優先権主張番号	10/014, 180		弁理士 柴田 沙希子
(32) 優先日	平成13年11月13日 (2001.11.13)	(72) 発明者	グッタ スリニヴァス ヴィ アール
(33) 優先権主張国	米国 (US)		オランダ国 5656 アーエー アイ
			ドーフエン プロフ ホルストラーン 6

最終頁に続く

(54) 【発明の名称】 アイテムの推薦器においてアイテムの近さを評価するための方法及び装置

(57) 【特許請求の範囲】

【請求項 1】

推薦器において用いられる方法であって、少なくとも1つのシンボリックなフィーチャを有しているアイテムを、1つ又は複数のグループに割り当て、前記推薦器のユーザのための視聴履歴を生成する方法において、前記推薦器が、

前記アイテムを第三者が利用した履歴を集めるステップと、

前記履歴における前記アイテムについて、前記アイテムと前記グループのそれぞれにおける少なくとも1つのアイテムとの対応するシンボリックなフィーチャ値間の距離を計算するステップであって、前記距離は、前記シンボリックなフィーチャ値の取り得る値それぞれに対応する、すべての実例のクラス分けの全体的な類似性に基づくものであるステップと

10

、
前記アイテムと前記グループのそれぞれにおける少なくとも1つのアイテムとの近さを決定する各前記フィーチャ値間の前記距離を総計するステップと、

前記アイテムを最小距離値に関連する前記グループに割り当てるステップと、

前記グループに基づいて、前記視聴履歴を生成するステップと、
を実施する、方法。

【請求項 2】

前記距離を計算する前記ステップは、前記シンボリックなフィーチャ間の前記距離を計算するために変更されたバリュエディファレンスメトリック技法を用いる、請求項 1 に記載の方法。

20

【請求項 3】

特定のシンボリックなフィーチャに関する 2 つの値 V 1 と V 2 との間の距離 は、

$$(V1, V2) = |C1_i / C1 - C2_i / C2|$$

によって与えられ、ここで、C 1 i は、V 1 がクラス i に分類された回数であり、C 1 は、V 1 がデータセットに出現した合計回数であり、C 2 i は、V 2 がクラス i に分類された回数であり、C 2 は、V 2 がデータセットに出現した合計回数である、請求項 1 に記載の方法。

【請求項 4】

前記アイテムは番組であり、クラス i は「視聴」クラス及び「非視聴」クラスであり、特定のシンボリックなフィーチャに関する 2 つの値 V 1 と V 2 との間の前記距離 は、

【数 1】

$$\delta(V1, V2) = \left| \frac{C1_watched}{C1_total} - \frac{C2_watched}{C2_total} \right| + \left| \frac{C1_not_watched}{C1_total} - \frac{C2_not_watched}{C2_total} \right|$$

によって与えられ、ここで、C 1 _w a t c h e d は V 1 が「視聴」クラスに分類された回数、C 1 _n o t _w a t c h e d は V 1 が「非視聴クラス」に分類された回数、C 1 _t o t a l は V 1 がデータセットに出現した合計回数であり、C 2 _w a t c h e d は V 2 が「視聴」クラスに分類された回数、C 2 _n o t _w a t c h e d は V 2 が「非視聴クラス」に分類された回数、C 2 _t o t a l は、V 2 がデータセットに出現した合計回数である、請求項 1 に記載の方法。

【請求項 5】

前記グループのそれぞれにおける前記少なくとも 1 つのアイテムは、予め計算されたグループの平均的なアイテムである、請求項 1 に記載の方法。

【請求項 6】

前記アイテムは番組である、請求項 1 に記載の方法。

【請求項 7】

前記アイテムはコンテンツである、請求項 1 に記載の方法。

【請求項 8】

前記アイテムは製品である、請求項 1 に記載の方法。

【請求項 9】

推薦器において用いられるシステムであって、少なくとも 1 つのシンボリックなフィーチャを有しているアイテムを、1 つ又は複数のグループに割り当て、前記推薦器のユーザのための視聴履歴を生成するシステムにおいて、前記推薦器が、

前記アイテムを第三者が利用した履歴を集める手段と、

前記履歴における前記アイテムについて、前記アイテムと前記グループのそれぞれにおける少なくとも 1 つのアイテムとの対応するシンボリックなフィーチャ値間の距離を計算する手段であって、前記距離は、前記シンボリックなフィーチャ値の取り得る値それぞれに対応する、すべての実例のクラス分けの全体的な類似性に基づくものである手段と、

前記アイテムと前記グループのそれぞれにおける少なくとも 1 つのアイテムとの近さを決定する各前記フィーチャ値の前記距離を総計する手段と、

前記アイテムを最小距離値に関連する前記グループに割り当てる手段と、

前記グループに基づいて、前記視聴履歴を生成する手段と、
を有する、システム。

【請求項 10】

10

20

30

40

50

推薦器において用いられるシステムであって、少なくとも1つのシンボリックなフィーチャを有しているアイテムを、1つ又は複数のグループに割り当て、前記推薦器のユーザのための視聴履歴を生成するシステムにおいて、コンピュータ読み取り可能なコードを記憶するためのメモリと、前記メモリに機能的に結合されたプロセッサとを有しており、前記プロセッサが、

前記アイテムを第三者が利用した履歴を集め、

前記履歴における前記アイテムについて、前記アイテムと前記グループのそれぞれにおける少なくとも1つのアイテムとの対応するシンボリックなフィーチャ値間の距離を、前記シンボリックなフィーチャ値の取り得る値それぞれに対応する、すべての実例のクラス分けの全体的な類似性に基づいて計算し、

前記アイテムと前記グループのそれぞれにおける少なくとも1つのアイテムとの近さを決定する各前記フィーチャ値間の前記距離を総計し、

前記アイテムを最小距離値に関連する前記グループに割り当て、

前記グループに基づいて、前記視聴履歴を生成するシステム。

【請求項11】

プログラム可能な装置がコンピュータプログラムを実行するときに請求項9に記載のシステムとして機能することを可能にするコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、「Method and Apparatus for Partitioning a Plurality of Items into Groups of Similar Items in a Recommender of Such Items」というタイトルの米国特許出願（代理人整理番号US010568）、「Method and Apparatus for Generating A Stereotypical Profile for Recommending Items of Interest Using Item-Based Clustering」というタイトルの米国特許出願（代理人整理番号US010569）、「Method and Apparatus for Recommending Items of Interest Based on Preferences of a Selected Third Party」というタイトルの米国特許出願（代理人整理番号US010572）、「Method and Apparatus for Recommending Items of Interest Based on Stereotype Preferences of Third Parties」というタイトルの米国特許出願（代理人整理番号US010575）及び「Method and Apparatus for Generating a Stereotypical Profile for Recommending Items of Interest Using Feature-Based Clustering」というタイトルの米国特許出願（代理人整理番号US010576）に関連し、これらの各出願は、本願と同時に
出願され、本発明の譲受人に譲渡され、参照によってここに盛り込まれる。

【0002】

本発明は、例えばテレビ番組のような興味のあるアイテムを推薦するための方法及び装置に関し、より具体的には、ユーザの購入又は視聴履歴が利用できるようになる前に興味のある番組及び他のアイテムを推薦するための技法に関する。

【背景技術】

【0003】

テレビの視聴者に利用可能なチャンネルの数は、このようなチャンネル上で利用可能な番組コンテンツの多様性ととも増加しているため、テレビの視聴者が興味のあるテレビ番組を識別することはますます骨の折れることになっている。電子番組ガイド（EPG）は、例えばタイトル、時間、日付及びチャンネルによって利用可能なテレビ番組を識別するとともに、個人化された好みに従って利用可能なテレビ番組がサーチされ又は分類されることを許すことによって興味のある番組の識別を容易にする。

【0004】

興味のあるテレビ番組及びその他のアイテムを推薦するための多数の推薦ツールが提案され又は提言されている。テレビ番組推薦ツールは、例えば、特定の視聴者にとって興味のある推薦された番組の組を得るために視聴者の好みを電子番組ガイドに適用する。

10

20

30

40

50

概して、テレビ番組推薦ツールは、黙示的若しくは明示的な技法を使用して又は前述のなんらかの組み合わせを使用して視聴者の好みを取得する。黙示的なテレビ番組推薦ツールは、視聴者の視聴履歴から導き出される情報に基づいて目立たないやり方でテレビ番組推薦を生成する。他方、明示的なテレビ番組推薦ツールは、視聴者プロフィールを導き出し推薦を生成するために、例えばタイトル、ジャンル、俳優、チャンネル及び日付/時間のような番組属性に関する好みについて明示的に視聴者に尋ねる。

【発明の開示】

【発明が解決しようとする課題】

【0005】

今日利用可能な推薦ツールは、ユーザが興味のあるアイテムを識別するのを助けるが、それらは多くの制約を受け、これらの制約は、克服される場合にはこのような推薦ツールの利便性及び性能を大幅に改善することができる。例えば、包括的に、明示的な推薦ツールは初期化するのが非常に長たらしく退屈であり、それぞれの新しいユーザは粒度の粗いレベルで彼らの好みを特定する非常に詳細な調査に耐える必要がある。黙示的なテレビ番組推薦ツールは、視聴の振る舞いを観察することによって目立たないようにプロフィールを導き出すが、それらは正確になるために長い時間を必要とする。加えて、このような黙示的なテレビ番組推薦ツールは、いかなる推薦をし始めるためにも視聴履歴の少なくとも最小量を必要とする。こうして、このような黙示的なテレビ番組推薦ツールは、推薦ツールが最初に取得されるときにはいかなる推薦もすることができない。

【0006】

従って、十分に個人化された視聴履歴が利用できるようになる前に目立たないように例えばテレビ番組のようなアイテムを推薦することができる方法及び装置が必要とされている。加えて、第三者の視聴履歴に基づいて所与のユーザのための番組推薦を生成するための方法及び装置が必要とされている。

【課題を解決するための手段】

【0007】

概して、例えばテレビ番組推薦のようなユーザにとって興味のあるアイテムを推薦するための方法及び装置が開示される。本発明の一つの見地によれば、例えばユーザが初めて推薦器を取得するときのように、ユーザの視聴履歴又は購入履歴が利用できるようになる前に推薦が生成される。最初に、1又は複数の第三者からの視聴履歴又は購入履歴が、特定のユーザにとって興味のあるアイテムを推薦するために使用される。

【0008】

第三者の視聴又は購入履歴は、代表的な視聴者により選択されたアイテムの一般的なパターンを反映する典型的(stereotype)なプロフィールを生成するために処理される。それぞれの典型的なプロフィールは、何らかの態様で互いに類似するアイテム(データポイント)のクラスタである。ユーザは、自分の興味に最も近いアイテムにより自分のプロフィールを初期化するために、興味のある1つ又は複数の典型を選択する。

【0009】

クラスタリングルーチンは、1つのクラスタにおけるポイント(例えばテレビ番組)が他のいかなるクラスタよりもそのクラスタの平均に近くなるように、第三者の視聴又は購入履歴(データセット)をクラスタに分ける。更に平均計算ルーチンが、クラスタのシンボリック(象徴的、symbolic)な平均を計算するために開示される。テレビ番組のような所与のデータポイントが、各クラスタの平均を使用して、データポイントと各クラスタとの間の距離に基づいてクラスタに割り当てられる。

【0010】

開示される距離計算ルーチンは、所与のテレビ番組と所与のクラスタの平均との間の距離に基づいて、テレビ番組の各クラスタに対する近さを評価する。計算された距離メトリックは、クラスタの範囲を決定するためにサンプルデータセットにおけるさまざまな標本間の区別を定量化する。バリュエディファレンスメトリック(value difference metric、VDM)技法又はその変形したものが、2つのテレビ番組間のフィーチャ値の間の距離

10

20

30

40

50

を計算するために用いられる。既知の変更されたVDM(modified VDM、MVDM)技法によれば、特定のフィーチャに関する2つの値の間の距離は、下式によって与えられる。

$$(V1, V2) = |C1i / C1 - C2i / C2|$$

ここで、V1及びV2は、考慮中のフィーチャに関する2つの可能な値である。例示的な実施例の番組推薦環境において、興味のあるクラスは、「視聴(watched)」及び「非視聴(not-watched)」である。概して、開示される距離計算ルーチンは、すべての分類について同じ相対頻度で値が生じる場合、それらの値を類似するものと識別する。こうして、2つの値は、それらがすべての可能な分類について類似する尤度(likelihoods)を与える場合は類似する。

10

【発明を実施するための最良の形態】

【0011】

本発明のより完全な理解及び本発明の他の特徴及び利点は、以下の詳細な説明及び図面を参照することにより得られる。

【0012】

図1は、本発明によるテレビ番組推薦器100を示す。図1に示すように、例示的なテレビ番組推薦器100は、特定の視聴者にとって興味のある番組を識別するために、図2に関連して以下に記述される番組データベース200における番組を評価する。推薦される番組の組は、例えば良く知られたオンスクリーン表示技法を使用するセットトップ端末/テレビジョン(図示せず)を使用して視聴者に対して提示されることができる。本発明は、ここでテレビ番組推薦のコンテキストにおいて説明されているが、本発明は、例えば視聴履歴又は購入履歴のようなユーザの振る舞いの評価に基づきいかなる自動的に生成される推薦にも適用されることができる。

20

【0013】

本発明の1つの特徴によれば、テレビ番組推薦器100は、例えばユーザが最初にテレビ番組推薦器100を取得するときのように、ユーザの視聴履歴140が利用できるようになる前に、テレビ番組推薦を生成することができる。図1に示すように、テレビ番組推薦器100は、まず、特定のユーザにとって興味のある番組を推薦するために1又は複数の第三者からの視聴履歴130を用いる。概して、第三者の視聴履歴130は、比較的大きい母集団を代表する例えば年齢、収入、性別及び教育のようなデモグラフィックスを有する1つ又は複数のサンプル母集団の視聴習慣に基づく。

30

【0014】

図1に示すように、第三者の視聴履歴130は、所与の母集団によって視聴される番組及び視聴されない番組の組からなる。視聴される番組の組は、所与の母集団によって実際に視聴される番組を観察することによって得られる。視聴されない番組の組は、例えば番組データベース200における番組をランダムに抽出することによって得られる。他の変形例において、視聴されない番組の組は、「An Adaptive Sampling Technique for Selecting Negative Examples for Artificial Intelligence Applications」というタイトルの米国特許出願第09/819,286号(2001年3月28日出願)の教示に従って得られる。前述の米国特許出願は本発明の譲受人に譲渡され、その内容はここに参照によって盛り込まれる。

40

【0015】

本発明の別の特徴によれば、テレビ番組推薦器100は、代表的な視聴者によって視聴されるテレビ番組の一般的なパターンを反映する典型的なプロファイルを生成するために第三者の視聴履歴130を処理する。以下に更に記述されるように、典型的なプロファイルは、何らかの態様で互いに類似するテレビ番組(データポイント)のクラスタである。こうして、所与のクラスタは、特定のパターンを示す第三者の視聴履歴130からのテレビ番組の特定のセグメントに対応する。

【0016】

第三者の視聴履歴130は、本発明に従って、ある特定のパターンを示す番組のクラス

50

タを提供するために処理される。そののち、ユーザは、最も関連する1つ又は複数の典型を選択することができ、それによって、自分の興味に最も近い番組により自分のプロフィールを初期化することができる。典型的なプロフィールは、それらの記録パターンに依存して、それぞれの個々のユーザの特定の個人の視聴振る舞いに対して適応し、発展し、フィードバックが番組に与えられる。一実施例において、番組スコアを決定する際、ユーザ自身の視聴履歴140からの番組には、第三者の視聴履歴130からの番組よりも高い重みを与えることができる。

【0017】

テレビ番組推薦器100は、中央処理ユニット(CPU)のようなプロセッサ115並びにRAM及び/又はROMのようなメモリ120を有する、パーソナルコンピュータ又はワークステーションのようないかなるコンピューティング装置としても具体化されることができる。テレビ番組推薦器100はまた、例えばセットトップ端末又はディスプレイ(図示せず)における特定用途向け集積回路(ASIC)として具体化されることもできる。加えて、テレビ番組推薦器100は、例えばカリフォルニア州サニーベールのTivo社から市販されているTivo(R)システム並びに「Method and Apparatus for Recommending Television Programming Using Decision Trees」というタイトルの米国特許出願第09/466,406号(1999年12月17日出願)、「Bayesian TV Show Recommender」というタイトルの米国特許出願第09/498,271号(2000年2月4日出願)、「Three-Way Media Recommendation Method and System」というタイトルの米国特許出願第09/627,139号(2000年7月27日出願)に添付の明細書に記載されているテレビ番組推薦器又はそれらの組み合わせのような、いかなる利用可能なテレビ番組推薦器としても具体化されることができる。前述の各々の米国特許出願の内容は本発明の特徴及び機能を実行するために変形されて参照によってここに盛り込まれる。

【0018】

図1に示され、図2乃至図8に関連して以下に更に記述されるように、テレビ番組推薦器100は、番組データベース200、典型的プロフィールプロセス300、クラスタリングルーチン400、平均計算ルーチン500、距離計算ルーチン600及びクラスタ性能評価ルーチン800を含む。概して、番組データベース200は、よく知られた電子番組ガイドとして具体化されることができ、所与の時間間隔において利用可能なそれぞれの番組に関する情報を記録する。典型的プロフィールプロセス300は、(i)代表的な視聴者によって視聴されるテレビ番組の一般的なパターンを反映する典型的なプロフィールを生成するために第三者の視聴履歴130を処理し、(ii)ユーザが最も関連する1つ又は複数の典型を選択し、それによって自分のプロフィールを初期化することを可能にし、(iii)選択された典型に基づいて推薦を生成する。

【0019】

クラスタリングルーチン400は、1つのクラスタにおけるポイント(テレビ番組)が他のいかなるクラスタよりもそのクラスタの平均(重心)に近くなるように第三者視聴履歴130(データセット)をクラスタに分けるために、典型的プロフィールプロセス300によって呼び出される。クラスタリングルーチン400は、クラスタのシンボリックな平均を計算するために平均計算ルーチン500を呼び出す。距離計算ルーチン600は、所与のテレビ番組と所与のクラスタの平均と間の距離に基づいてそれぞれのクラスタに対するテレビ番組の近さを評価するために、クラスタリングルーチン400によって呼び出される。最後に、クラスタリングルーチン400は、クラスタを生成するためのストップ基準が満たされたときを決定するために、クラスタリング性能評価ルーチン800を呼び出す。

【0020】

図2は、図1の番組データベース(電子番組ガイド)200からのサンプルテーブルである。上述したように、番組データベース200は、所与の時間間隔において利用可能であるそれぞれの番組ごとの情報を記録する。図2に示すように、番組データベース200は、例えばレコード205乃至220のような複数のレコードを含み、それぞれのレコー

10

20

30

40

50

ドは、所与の番組に関連する。それぞれの番組について、番組データベース200は、フィールド240及び245においてそれぞれ番組に関連する日付/時間及びチャンネルを示す。加えて、それぞれの番組に関するタイトル、ジャンル及び俳優は、フィールド250、255及び270においてそれぞれ識別される。番組の持続期間及び説明のようによく知られた付加的なフィーチャ(図示せず)もまた、番組データベース200に含められることができる。

【0021】

図3は、本発明の特徴を取り入れる典型的プロファイルプロセス300の例示的な実現例を示すフローチャートである。上述したように、典型的プロファイルプロセス300は、(i)代表的な視聴者によって視聴されるテレビ番組の一般的なパターンを反映する典型的なプロファイルを生成するために第三者の視聴履歴130を処理し、(ii)ユーザが最も関連する1つ又は複数の典型を選択し、それによって自分のプロファイルを初期化することを可能にし、(iii)選択された典型に基づいて推薦を生成する。第三者の視聴履歴130の処理は、オフラインで実施されてもよく、例えば工場で実施されてもよく、テレビ番組推薦器100は、ユーザによる選択のために生成された典型的なプロファイルをインストールされて、ユーザに提供されることができるとに留意する。

10

【0022】

こうして、図3に示すように、典型的プロファイルプロセス300は、まずステップ310の間に、第三者の視聴履歴130を集める。そののち、ステップ320の間に、典型的プロファイルプロセス300は、典型的なプロファイルに対応する番組のクラスタを生成するために、図4に関連して以下に記述されるクラスタリングルーチン400を実行する。以下に更に記述されるように、例示的なクラスタリングルーチン400は、視聴履歴データセット130に対して、例えば「k-means(k平均)」クラスタルーチンのような非監督型データクラスタリングアルゴリズムを用いることができる。上述したように、クラスタリングルーチン400は、1つのクラスタにおけるポイント(テレビ番組)が他のいかなるクラスタよりもそのクラスタの平均(重心)に近くなるように、第三者視聴履歴130(データセット)をクラスタに分ける。

20

【0023】

ステップ330の間、典型的プロファイルプロセス300は、それぞれの典型的なプロファイルの特徴付ける1つ又は複数のラベルをそれぞれのクラスタに割り当てる。1つの例示的な実施例において、クラスタの平均は、クラスタ全体について代表的なテレビ番組になり、平均番組のフィーチャは、クラスタにラベルをつけるために使用されることができる。例えば、テレビ番組推薦器100は、ジャンルがそれぞれのクラスタについて優勢な又は規定するフィーチャであるように構成されることができるとに留意する。

30

【0024】

ステップ340の間、ラベルをつけられた典型的なプロファイルは、ユーザの興味に最も近い1つ又は複数の典型的なプロファイルの選択のために、それぞれのユーザに対して提示される。それぞれの選択されたクラスタを構成する番組は、当該典型の「一般的な視聴履歴」と考えられることができ、それぞれのクラスタについて典型的なプロファイルを構築するために使用されることができるとに留意する。こうして、ステップ350の間、選択された典型的なプロファイルからの番組からなる視聴履歴がユーザのために生成される。最後に、ステップ360の間、前のステップで生成された視聴履歴が、番組推薦を得るために番組推薦器に適用される。番組推薦器は、当業者にとって明らかであるように、上述した番組推薦器のようないかなる通常の番組推薦器としてここで変形されて具体化されてもよい。プログラム制御は、ステップ370の間に終了する。

40

【0025】

図4は、本発明の特徴を取り入れるクラスタリングルーチン400の例示的な実現例を示すフローチャートである。上述したように、ステップ320の間、クラスタリングルーチン400は、1つのクラスタにおけるポイント(テレビ番組)が他のいかなるクラスタよりもそのクラスタの平均(重心)に近くなるように第三者視聴履歴130(データセッ

50

ト)をクラスタに分けるために、典型的プロファイルプロセス300によって呼び出される。概して、クラスタリングルーチンは、サンプルデータセットにおける標本のグループを見つける監督されないタスクに注目する。本発明は、k-meansクラスタリングアルゴリズムを使用して、データセットをk個のクラスタに分ける。以下に記述されるように、クラスタリングルーチン400への2つの主なパラメータは、(i)図6に関連して以下に記述される最も近いクラスタを見つけるための距離メトリック及び(ii)生成すべきクラスタの数kである。

【0026】

標本データの更なるクラスタリングが分類の精度のいかなる改善も与えないときに安定したkに達するという条件で、例示的なクラスタリングルーチン400は、kの動的な値を使用する。加えて、クラスタサイズは、空のクラスタが記録されるところまでインクリメントされる。こうして、クラスタの自然なレベルに達するとき、クラスタリングは止まる。

10

【0027】

図4に示すように、クラスタリングルーチン400は、まずステップ410の間に、k個のクラスタを確立する。例示的なクラスタリングルーチン400は、クラスタの最小数、例えば2を選択することから始まる。この固定の数について、クラスタリングルーチン400は、視聴履歴データセット130全体を処理し、いくつかの繰り返しを通して、安定していると考えられうる2つのクラスタに達する(すなわちアルゴリズムが別の繰り返しを通るとしても、いかなる番組もあるクラスタから別のクラスタに移動しない)。ステップ420の間、現在のk個のクラスタは、1つ又は複数の番組により初期化される。

20

【0028】

1つの例示的な実現例において、ステップ420の間、クラスタは、第三者視聴履歴130から選択されるシード番組により初期化される。クラスタを初期化するための番組は、ランダムに又は逐次的に選択されることができる。逐次的な実現例において、クラスタは、視聴履歴130における1番目の番組から始まる番組を用いて又は視聴履歴130におけるランダムなポイントから始まる番組を用いて初期化されることができる。更に他の変形例において、それぞれのクラスタを初期化する番組の数が変えられてもよい。最後に、クラスタは、第三者視聴履歴130における番組からランダムに選択されたフィーチャ値からなる1つ又は複数の「仮定的な(hypothetical)」番組を用いて初期化されてもよい。

30

【0029】

そののち、ステップ430の間、クラスタリングルーチン400は、それぞれのクラスタの現在の平均を計算するために、図5に関連して後述される平均計算ルーチン500を始める。ステップ440の間、クラスタリングルーチン400は、第三者視聴履歴130におけるそれぞれの番組の各クラスタに対する距離を決定するために、図6に関連して記述される距離計算ルーチン600を実行する。ステップ460の間、視聴履歴130におけるそれぞれの番組は、最も近いクラスタに割り当てられる。

【0030】

ステップ470の間、番組があるクラスタから別のクラスタへ移動したかどうか決定するためのテストが実施される。番組があるクラスタから別のクラスタに移動したとステップ470の間に決定される場合、プログラム制御はステップ430に戻り、クラスタの安定した組が識別されるまで上述した態様で続く。しかしながら、いかなる番組もあるクラスタから別のクラスタにステップ470の間に移動しなかったと決定される場合、プログラム制御はステップ480に進む。

40

【0031】

ステップ480の間、特定された性能基準が満たされたか又は空のクラスタが識別されるか(総称して「ストップ基準」と呼ぶ)を決定するための他のテストが実施される。ステップ480の間にストップ基準が満たされていないと決定される場合、ステップ485の間にkの値がインクリメントされ、プログラム制御はステップ420に戻り、上述した

50

態様で続く。しかしながら、ステップ480の間にストップ基準が満たされたと決定される場合、プログラム制御が終了する。ストップ基準の評価は、図8に関連して以下に更に記述される。

【0032】

例示的なクラスタリングルーチン400は、ただ1つのクラスタに番組を配置し、それゆえクリスタルクラスタと呼ばれるものを生成する。他の変形例は、特定の標本(テレビ番組)が部分的に多くのクラスタに属することを可能にするファジークラスタリングを用いる。ファジークラスタリング方法において、テレビ番組は、テレビ番組がクラスタ平均にどれくらい近いかを表す重みを割り当てられる。重みは、クラスタ平均からのテレビ番組の距離の逆二乗に依存する。単一のテレビ番組に関連するすべてのクラスタ重みの合計は、100%にならなければならない。

10

【0033】

クラスタのシンボリックな平均の計算

図5は、本発明の特徴を取り入れる平均計算ルーチン500の例示的な実現例を示すフローチャートである。上述したように、平均計算ルーチン500は、クラスタのシンボリックな平均を計算するために、クラスタリングルーチン400によって呼び出される。数値データについて、平均は、分散を最小にする値である。シンボリックなデータに概念を拡張するとき、クラスタの平均は、クラスタ内の分散(それゆえクラスタの半径又は範囲)を最小にする x_{μ} の値を求めることによって規定されることができる。

【数2】

20

$$\text{Var}(J) = \sum_{i \in J} (x_i - x_{\mu})^2 \quad (1)$$

$$\text{Cluster radius } R(J) = \sqrt{\text{Var}(J)} \quad (2)$$

ここで、Jは、同じクラス(視聴クラス又は非視聴クラス)からのテレビ番組のクラスタであり、 x_i は、ショー(番組)iについてのシンボリックなフィーチャ値であり、 x_{μ} は、 $\text{Var}(J)$ を最小にするようなJにおけるテレビ番組の1つからのフィーチャ値である。

30

【0034】

こうして、図5に示すように、平均計算ルーチン500は、まずステップ510の間、所与のクラスタJにおける現在の番組を識別する。ステップ520の間、考慮中の現在のシンボリックな属性について、クラスタJの分散が、それぞれの可能なシンボリックな値 x_{μ} について方程式(1)を使用して計算される。ステップ530の間、分散を最小にするシンボリックな値 x_{μ} が平均値として選択される。

【0035】

ステップ540の間、考慮されるべき他のシンボリックな属性があるかを決定するためのテストが実施される。考慮されるべき他のシンボリックな属性があるとステップ540の間に決定される場合、プログラム制御はステップ520に戻り、上述した態様で続く。しかしながら、考慮されるべき他のシンボリックな属性がないとステップ540の間に決定される場合、プログラム制御はクラスタリングルーチン400に戻る。

40

【0036】

計算的に、Jにおけるそれぞれのシンボリックなフィーチャ値が x_{μ} として試され、分散を最小にするシンボリックな値がクラスタJにおける考慮中のシンボリックな属性の平均になる。可能である平均計算には2つのタイプがあり、すなわちショー(番組)に基づく平均とフィーチャに基づく平均とがある。

【0037】

フィーチャに基づくシンボリックな平均

50

ここに記述される例示的な平均計算ルーチン500はフィーチャに基づく。ここで、シンボリックな属性に関する平均は、その可能な値のうちの一つでなければならないので、結果として得られるクラスタ平均は、クラスタJにおける標本(番組)から導き出されるフィーチャ値からなる。しかしながら、クラスタ平均は、「仮定的な」テレビ番組でありうる点に留意することが重要である。この仮定的な番組のフィーチャ値は、標本のうちの一つから引き出されるチャンネル値(例えばEBC)及び標本のうちの別のものから引き出されるタイトル値(例えば実際にEBCで放送されることはないBBCワールドニュース)を含むことができる。このように、最小分散を呈するいかなるフィーチャ値も、当該フィーチャの平均を表すために選択される。ステップ540の間にすべてのフィーチャ(すなわちシンボリックな属性)が考慮されたと決定されるまで、平均計算ルーチン500はすべてのフィーチャ位置について繰り返される。こうして結果として得られた仮定的な番組は、クラスタの平均を表すために使用される。

10

【0038】

番組に基づくシンボリックな平均

他の変形例では、分散に関する方程式(1)において、 x_i は、テレビ番組i自体でありえ、同様に、 x_j は、クラスタJにおける番組の組にわたる分散を最小にするクラスタJにおける1つ又は複数の番組である。この例では、個々のフィーチャ値でなく番組間の距離が、最小にされるべき関連するメトリックである。加えて、この例において結果として得られる平均は、仮定的な番組でなく、組Jからまさに選ばれた番組である。クラスタJにおけるすべての番組にわたる分散を最小にする、こうして見つけられたクラスタJにおけるいかなる番組も、クラスタの平均を表すために使用される。

20

【0039】

複数の番組を使用するシンボリックな平均

上述の例示的な平均計算ルーチン500は、(フィーチャ又は番組に基づく実現例のいずれにせよ)それぞれの可能なフィーチャについて単一のフィーチャ値を使用してクラスタの平均を特徴付ける。しかしながら、平均は、もはや当該クラスタの代表的なクラスタ中心ではないので、平均計算の間にそれぞれのフィーチャについてただ1つのフィーチャ値に依存することは、多くの場合、不適当なクラスタリングをもたらす。言い換えると、ただ1つの番組によって、クラスタを表すことは望ましくないことがあり、むしろ、1つ又は複数の平均を表す複数の番組が、クラスタを表すために使用されることができ。こうして、他の変形例において、クラスタは、それぞれの可能なフィーチャに関する複数の平均又は複数のフィーチャ値によって表現されることができ。こうして、分散を最小にするN個のフィーチャ(フィーチャに基づくシンボリックな平均の場合)又はN個の番組(番組に基づくシンボリックな平均の場合)が、ステップ530の間に選択される。ここで、Nは、クラスタの平均を表すために使用される番組の数である。

30

【0040】

番組とクラスタと間の距離計算

前述のように、所与のテレビ番組と所与のクラスタの平均と間の距離に基づいて、それぞれのクラスタに対するテレビ番組の近さを評価するために、距離計算ルーチン600がクラスタリングルーチン400によって呼び出される。計算された距離メトリックは、クラスタの範囲を決定するためにサンプルデータセットにおけるさまざまな標本間の区別を定量化する。ユーザプロファイルをクラスタ分けすることができるようにするために、視聴履歴における任意の2つのテレビ番組間の距離が計算されなければならない。概して、互いに近いテレビ番組は、1つのクラスタに入る傾向がある。数値的に表されるベクトル間の距離を計算するために多くのかなり直接的な技法が存在し、例えばユークリッド距離、マンハッタン距離及びマハラノビス距離がある。

40

【0041】

しかしながら、テレビ番組は主としてシンボリックなフィーチャ値からなるので、既存の距離計算技法はテレビ番組ベクトルの場合には使用されることができない。例えば、2001年3月22日午後8時にEBCで放送される「Friends」の一話及び2001年

50

3月25日午後8時にF E Xで放送される「The Simons」の一話のような2つのテレビ番組は、以下のフィーチャベクトルを使用して表現されることができる。

タイトル：F i e n d s	タイトル：S i m o n s
チャンネル：E B C	チャンネル：F E X
放送日：2 0 0 1 - 0 3 - 2 2	放送日：2 0 0 1 - 0 3 - 2 5
放送時間：2 0 0 0	放送時間：2 0 0 0

【 0 0 4 2 】

明らかに、既知の数値距離メトリックは、フィーチャ値「E B C」と「F E X」との間の距離を計算するために使用されることができない。バリュウディファレンスマトリック (V D M) は、シンボリックなフィーチャで表されるドメインにおけるフィーチャの値の間の距離を測るための既存の技法である。V D M技法は、それぞれのフィーチャのそれぞれの取り得る値についてすべての実例の分類 (classification)の全体的な類似性を考慮する。この方法を使用して、フィーチャのすべての値の間の距離を規定するマトリクスが、トレーニングセットにおける標本に基づいて統計的に導き出される。シンボリックなフィーチャ値の間の距離を計算するためのV D M技法に関するより詳細な説明のため、例えばSt anfill及びWaltz、「Toward Memory-Based Reasoning」、Communications of the ACM、2 9:12、1213-1228、1986年、を参照されたい。この内容は、参照によってここに盛り込まれる。

10

【 0 0 4 3 】

本発明は、2つのテレビ番組間又は興味のある他のアイテム間のフィーチャ値の間の距離を計算するためにV D M技法又はその変形を用いる。本来のV D M提案は、2つのフィーチャ値の間の距離計算において重みタームを使用し、これは、距離メトリックを非対照にする。変更されたV D M (M V D M) は、距離マトリックスを対称にするために重みタームを省く。シンボリックなフィーチャ値の間の距離を計算するためのM V D M技法のより詳細な説明のため、Cost及びSalzberg、「A Weighted Nearest Neighbor Algorithm For Learning With Symbolic Features」、Machine Learning、Vol. 10、57-58、ボストン、マサチューセッツ、Kluwer Publishers、1993年、を参照されたい。この内容は、参照によってここに盛り込まれる。

20

【 0 0 4 4 】

M V D Mに従って、特定のフィーチャに関する2つの値V 1とV 2と間の距離が、下式によって与えられる。

30

$$(V 1, V 2) = | C 1 i / C 1 - C 2 i / C 2 | \quad (3)$$

【 0 0 4 5 】

本発明の番組推薦環境において、M V D M方程式(3)は、特に「視聴」及び「非視聴」クラスを扱うために変形される。

【数3】

$$\delta(V1, V2) = \left| \frac{C1_watched}{C1_total} - \frac{C2_watched}{C2_total} \right| + \left| \frac{C1_not_watched}{C1_total} - \frac{C2_not_watched}{C2_total} \right| \quad (4)$$

40

【 0 0 4 6 】

方程式(4)において、V 1及びV 2は、考慮中のフィーチャに関する2つの可能な値である。上述の例を続けると、フィーチャ「チャンネル」に関して、第1の値V 1は「E B C」に等しく、第2の値V 2は「F E X」に等しい。これらの値の間の距離は、標本が分類されるすべてのクラスにわたる合計である。本発明の例示的な番組推薦器の実施例の場合の関連するクラスは、「視聴」クラス及び「非視聴」クラスである。C 1 iは、V 1 (

50

EBC) がクラス i に分類された回数であり (1 に等しい i は、視聴クラスを意味する)、 $C1$ ($C1_total$) は、 $V1$ がデータセットにおいて出現した合計回数である。値「 r 」は、定数であり、通常 1 にセットされる。

【0047】

方程式 (4) によって規定されるメトリックは、値がすべての分類について同じ相対頻度で生じる場合、それらの値を類似するものとして識別する。ターム $C1_i / C1$ は、当該フィーチャが値 $V1$ を有すると仮定して、その中央残余 (central residue) が i として分類されるであろう尤度を表現する。こうして、2 つの値は、それらがすべての可能な分類について類似の尤度を与える場合には類似する。方程式 (4) は、すべての分類にわたるこれらの尤度の差の合計を求めることによって、2 つの値の間の全体的な類似性を計算する。2 つのテレビ番組間の距離は、2 つのテレビ番組ベクトルの対応するフィーチャ値間の距離の合計である。

10

【0048】

図 7A は、フィーチャ「チャンネル」に関連するフィーチャ値に関する距離テーブルの一部である。図 7A は、それぞれのクラスについてそれぞれのチャンネルフィーチャ値の出現の数を含む。図 7A に示される値は、例示的な第三者視聴履歴 130 から取り出されたものである。

【0049】

図 7B は、MVD M 方程式 (4) を使用して図 7A に示される例示的なカウントから計算されるそれぞれのフィーチャ値対の間の距離を示す。直観的に、EBC 及び ABS は、それらが大部分は視聴クラスに出現し、非視聴クラスに出現しない (ABS は小さい非視聴成分を有する) ので、互いに「近い」とすべきである。図 7B は、EBC と ABS との間の小さい (ゼロでない) 距離によりこの直感的事実を確認する。他方、ASPN は、大部分は非視聴クラスに出現し、それゆえこのデータセットについては EBC 及び ABS のどちらに対しても「遠い」とすべきである。図 7B は、最大可能距離 2.0 のところ、EBC と ASPN との間の距離 1.895 を含む。同様に、ABS と ASPN と間の距離は高く、値 1.828 をもつ。

20

【0050】

こうして、図 6 に示すように、まずステップ 610 の間、距離計算ルーチン 600 は、第三者視聴履歴 130 における番組を識別する。ステップ 620 の間、距離計算ルーチン 600 は、考慮中の現在番組について、(平均計算ルーチン 500 により決定された) それぞれのクラス平均の対応するフィーチャに対するそれぞれのシンボリックなフィーチャ値の距離を計算するために方程式 (4) を使用する。

30

【0051】

ステップ 630 の間、現在番組とクラス平均と間の距離が、対応するフィーチャ値の間の距離を総計することによって計算される。ステップ 640 の間、考慮されるべき他の番組が第三者視聴履歴 130 にあるかどうか決定するためのテストが実施される。ステップ 640 の間に考慮されるべき他の番組が第三者視聴履歴 130 にあると決定される場合、次の番組がステップ 650 の間に識別され、プログラム制御はステップ 620 に進み、上述した態様で続く。

40

【0052】

しかしながら、ステップ 640 の間、考慮されるべき他の番組が第三者視聴履歴 130 にはないと決定される場合、プログラム制御はクラスタリングルーチン 400 へ戻る。

【0053】

「複数の番組から導き出されるシンボリックな平均」という題目のサブセクションにおいて上述したように、クラスタの平均は、(フィーチャ又は番組に基づく実現例のいずれにせよ) それぞれの可能なフィーチャについて複数のフィーチャ値を使用して特徴付けられることができる。複数の平均からの結果は、投票によるコンセンサス決定に達するために距離計算ルーチン 600 の変形によってプールされる。例えば、ステップ 620 の間、番組の所与のフィーチャ値と、さまざまな平均についての各々の対応するフィーチャ値との

50

間で、距離が計算される。最小距離結果はプールされ、例えばコンセンサス決定に達するために多数決又は混合エキスパート (mixture of experts) を用いることによって、投票のために使用される。このような技法に関するより詳細な説明のため、例えば、J.Kittler 他、「Combining Classifiers」、Proc. of the 13th Int'l Conf. on Pattern Recognition、Vol. 11、897-901、ウィーン、オーストリア、1996年、を参照されたい。その内容は参照によりここに盛り込まれる。

【0054】

ストップ基準

上述したように、クラスタリングルーチン400は、クラスタを生成するためのストップ基準が満たされたときを決定するために、図8に示されるクラスタリング性能評価ルーチン800を呼び出す。標本データの更なるクラスタリングが分類精度のいかなる改善も与えないときに安定したkに達するという条件で、例示的なクラスタリングルーチン400は、kの動的な値を用いる。加えて、クラスタサイズは、空のクラスタが記録されるところまでインクリメントされることができる。こうして、クラスタの自然なレベルに達するときに、クラスタリングがストップする。

10

【0055】

例示的なクラスタリング性能評価ルーチン800は、クラスタリングルーチン400の分類精度をテストするために、第三者視聴履歴130からの番組のサブセット(テストデータセット)を使用する。テストセットにおけるそれぞれの番組について、クラスタリング性能評価ルーチン800は、当該番組に最も近いクラスタ(そのクラスタ平均が最も近い)を決定し、そのクラスタに関するクラスラベルと考慮中の番組とを比較する。合致したクラスラベルのパーセンテージは、クラスタリングルーチン400の精度につながる。

20

【0056】

こうして、図8に示すようにクラスタリング性能評価ルーチン800は、まずステップ810の間に、テストデータセットの役目を果たすための番組のサブセットを第三者視聴履歴130から集める。そののち、ステップ820の間、クラスタにおける視聴される及び視聴されない番組のパーセンテージに基づいて、クラスラベルがそれぞれのクラスタに割り当てられる。例えばクラスタにおける大部分の番組が視聴されるものである場合、そのクラスタには「視聴」ラベルが割り当てられることができる。

30

【0057】

テストセットにおけるそれぞれの番組に最も近いクラスタが、ステップ830の間に識別され、割り当てられたクラスタについてのクラスラベルは、番組が実際に視聴されたか否かと比較される。複数の番組がクラスタの平均を表示するために使用される実現例において、(それぞれの番組に対する)平均的な距離又は投票スキームが用いられることができる。プログラム制御がクラスタリングルーチン400へ戻る前に、合致したクラスラベルのパーセンテージが、ステップ840の間に決定される。分類精度があらかじめ規定された閾値に達した場合、クラスタリングルーチン400が終了する。

【0058】

図示されると共にここに記述される実施例及び変形例は、単に本発明の原理を説明するだけであり、さまざまな変更が、本発明の範囲及び精神から逸脱することなく当業者により実現されることが理解されるべきである。

40

【図面の簡単な説明】

【0059】

【図1】本発明によるテレビ番組推薦器の概略ブロック図。

【図2】図1の例示的な番組データベースからのサンプルテーブルを示す図。

【図3】本発明の原理を具体化する図1の典型的プロファイルプロセスを示すフローチャート。

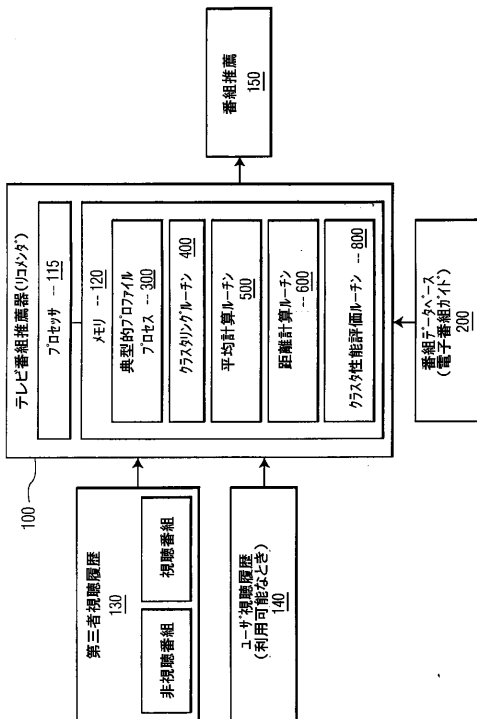
【図4】本発明の原理を具体化する図1のクラスタリングルーチンを示すフローチャート。

【図5】本発明の原理を具体化する図1の平均計算ルーチンを示すフローチャート。

50

【図 6】本発明の原理を具体化する図 1 の距離計算ルーチンを示すフローチャート。
 【図 7 A】それぞれのクラスについて、それぞれのチャンネルフィーチャ値の出現の数を示す例示的なチャンネルフィーチャ値出現テーブルからのサンプルテーブルを示す図。
 【図 7 B】図 7 A に示される例示的なカウントから計算されるそれぞれのフィーチャ値対の間の距離を示す例示的なフィーチャ値対距離テーブルからのサンプルテーブルを示す図。
 【図 8】本発明の原理を具体化する図 1 のクラスタリング性能評価ルーチンを示すフローチャート。

【図 1】



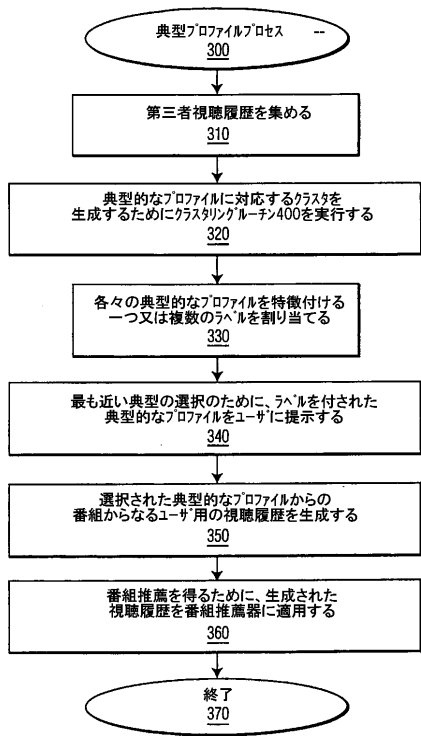
【図 2】

番組データベース - 200

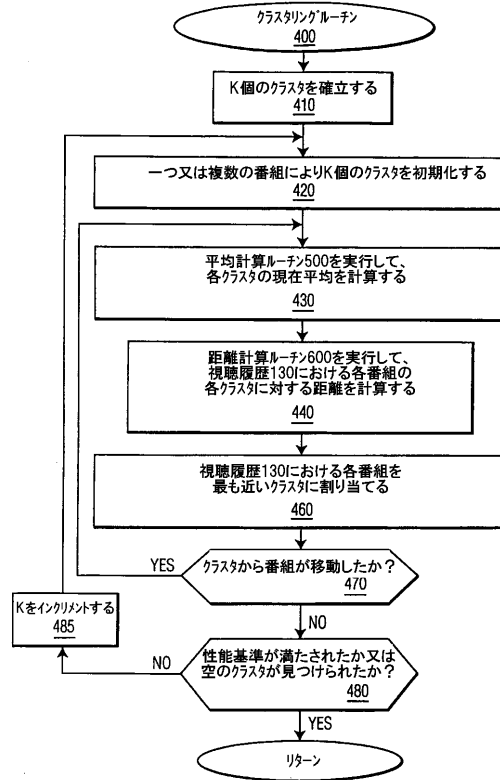
日付/時間	チャンネル	タイトル	ジャンル	俳優
240	245	250	255	270
11/18/99 -- 8:00 P.M.	CH1	LUCY	COMEDY	CLINT DENIRO
11/18/99 -- 8:30 P.M.	CH1	AL'S FAMILY	SITCOM	JENNIFER COX
...				
11/18/99 -- 9:00 P.M.	CH3	YOUR HOUSE	DRAMA	LUCY VANCE

205
210
220

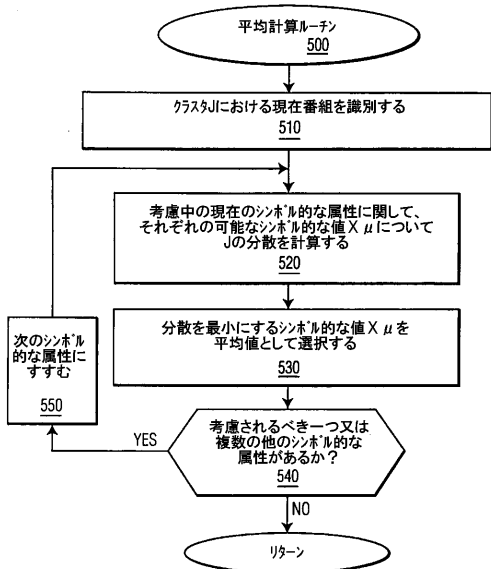
【図3】



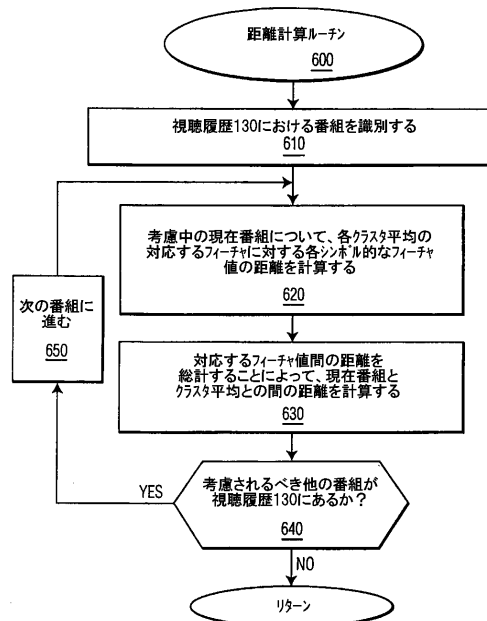
【図4】



【図5】



【図6】



【図7A】

チャンネルフィーチャ出現テーブル -- 700

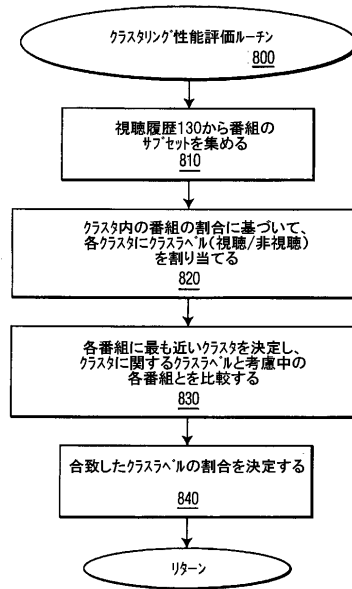
フィーチャ値	クラス	
	視聴	非視聴
EBC	353	0
ASPN	1	18
ABS	145	5

【図7B】

フィーチャ値対の距離のテーブル -- 750

	EBC	ASPN	ABS
EBC	0	1.895	0.066
ASPN	1.895	0	1.828
ABS	0.066	1.828	0

【図8】



フロントページの続き

(72)発明者 クラパティ カウシャル
オランダ国 5 6 5 6 アーアー アインドーフェン プロフ ホルストラーン 6

合議体

審判長 山崎 達也

審判官 原 秀人

審判官 田中 秀人

- (56)参考文献 特開2001-306612(JP,A)
特開2001-076002(JP,A)
特開2001-043233(JP,A)
特開2000-293531(JP,A)
河村晃好,外3名,「グループ嗜好モデルと視聴履歴を利用したコンテンツ検索サーバの試作」,
電子情報通信学会技術研究報告,日本,社団法人電子情報通信学会,平成13年7月11日,
第101巻 第192号,p.9-16
渡部勇,「緩い協調:協調情報フィルタリングシステム」,情報処理学会研究報告,日本,社団
法人情報処理学会,平成3年3月8日発行,第91巻 第18号,p.179-186
Scott Cost, Steven Salzberg, "A Weighted Near
est Neighbor Algorithm for Learning with Sy
mbolic Features", Machine Learning, 米国, Kluwer
Academic Publishers, 1993年1月, Volume 10 , Issu
e 1, Pages: 57-78

(58)調査した分野(Int.Cl., DB名)

G06F 17/30