

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4879178号
(P4879178)

(45) 発行日 平成24年2月22日(2012.2.22)

(24) 登録日 平成23年12月9日(2011.12.9)

(51) Int.Cl. F 1
G06T 7/00 (2006.01) G06T 7/00 350Z

請求項の数 28 (全 20 頁)

(21) 出願番号	特願2007-529054 (P2007-529054)	(73) 特許権者	504194580
(86) (22) 出願日	平成17年8月1日(2005.8.1)		石川 博
(65) 公表番号	特表2008-508645 (P2008-508645A)		埼玉県富士見市ふじみ野東1丁目22-2
(43) 公表日	平成20年3月21日(2008.3.21)		-102
(86) 国際出願番号	PCT/IB2005/052570	(72) 発明者	石川 博
(87) 国際公開番号	W02006/013549		愛知県東海市加木屋町北鹿持34-93
(87) 国際公開日	平成18年2月9日(2006.2.9)		
審査請求日	平成20年7月22日(2008.7.22)	審査官	新井 則和
(31) 優先権主張番号	60/592,911	(56) 参考文献	特開平04-276785 (JP, A)
(32) 優先日	平成16年8月2日(2004.8.2)	(58) 調査した分野(Int.Cl., DB名)	G06T 7/00
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 自動パターン解析のための方法と装置

(57) 【特許請求の範囲】

【請求項1】

記憶手段と処理手段を備えた情報処理システムにより実行されるパターン解析方法であって、

少なくとも1つの第1データを受信する第1工程と、

該第1データを、データ、写像、およびパターンの記録可能な前記記憶手段上のデータ構造に記録する第2工程と、

該第1データの構造と解釈に応じて少なくとも1つの原始写像を決定する第3工程と、

該原始写像を前記データ構造に記録する第4工程と、

前記データ構造に記録された少なくとも1つの第2データを選択する第5工程と、

該第2データから少なくとも1つの第3データを誘導する複数の手続きから少なくとも1つの手続きを選択する第6工程と、

該手続きに従い前記第2データから前記第3データを誘導する第7工程と、

前記データ構造に該第3データを記録する第8工程と

を有し、前記複数の手続きに、前記データ構造に記録された少なくとも1つの第1写像を所定の方法により選択し該第1写像を前記第2データに適用することにより前記第3データを誘導することを特徴とする第1手続きと、

前記第2データ内に表現された複数の集合の直積を少なくとも1つとり該直積を前記第3データ内に表現することにより該第3データを誘導することを特徴とする第2手続きと

を含むことを特徴とするパターン解析方法であって、

10

20

前記データ構造に記録された少なくとも1つの第4データ内に所定の方法により少なくとも1つの第1パターンを探す第9工程をさらに有することを特徴としたパターン解析方法。

【請求項2】

前記第9工程は前記第4データのエンтроピーに従って前記第1パターンを探すことを特徴とした請求項1記載のパターン解析方法。

【請求項3】

前記第1パターンが見つければ該第1パターンを前記データ構造に記録する第10工程と、

前記データ構造に記録された少なくとも1つの第2パターンをパターン解析結果として提供する第11工程

をさらに有することを特徴とした請求項1ないし2いずれか1項に記載のパターン解析方法。

【請求項4】

前記第1パターンが見つければ該第1パターンに対応する理想化されたデータである少なくとも1つの第5データを所定の方法により生成し前記データ構造に記録する第12工程をさらに有することを特徴とした請求項1ないし3いずれか1項に記載のパターン解析方法。

【請求項5】

前記第12工程は前記第4データ内に表現された少なくとも1つの第1確率測度を選択しエンтроピーのより低い第2確率測度を該第1確率測度から生成し該第2確率測度を前記第5データ内に表現するか、

前記第4データ内に表現された少なくとも1つの第3確率測度を選択し該第3確率測度を集中させて少なくとも1つの第4確率測度を生成し該第4確率測度を前記第5データ内に表現するか、

前記第4データ内に表現された少なくとも1つの第5確率測度を選択し該第5確率測度内の少なくとも1つの確率の集中に各々対応した複数の確率測度を生成し該複数の確率測度を前記第5データ内に表現するか、

前記第4データ内の近似的に繰り返すパターンをより正確に前記第5データ内に繰り返させるか

の少なくとも1つにより前記第5データを生成することを特徴とした請求項4記載のパターン解析方法。

【請求項6】

前記第1パターンが見つければ前記第4データに結び付けられた手続き及び写像に従ってパターン写像を決定し前記データ構造に記録する第13工程をさらに有し、

前記第8工程は選択された前記手続きと該手続きで写像が使用されていれば該写像を前記第3データと結び付けて前記データ構造に記録することを特徴とした請求項1ないし5いずれか1項に記載のパターン解析方法。

【請求項7】

一連の工程を所定の停止条件が満たされるまで繰り返す第14工程をさらに有し、該一連の工程は前記第5ないし14工程のうち前記第14工程以外の少なくとも1つを含むことを特徴とした請求項1ないし6いずれか1項に記載のパターン解析方法。

【請求項8】

前記複数の手続きに、

前記データ構造に記録された少なくとも1つの第2写像を選択し前記第2データ内に表現された少なくとも1つの第1集合の該第2写像による逆像をとり前記第3データ内に該逆像を表現することにより該第3データを誘導することを特徴とする第3手続きをさらに含むことを特徴とした請求項1ないし7いずれか1項に記載のパターン解析方法。

【請求項9】

前記複数の手続きに、

10

20

30

40

50

前記第 2 データ内に表現された少なくとも 1 つの第 2 集合の少なくとも 1 つの部分集合をとり前記第 3 データ内に該部分集合を表現することにより該第 3 データを誘導することを特徴とする第 4 手続きをさらに含むことを特徴とした請求項 1 ないし 8 いずれか 1 項に記載のパターン解析方法。

【請求項 10】

前記原始写像が、

恒等写像、定数写像、等号写像、積写像、複数の写像の積写像を与える写像、引戻し演算写像、射影写像、対角写像、置換写像、写像合成写像、評価写像、複数の低位写像を組み合わせ高位写像を与える写像、CURRY写像、論理演算写像、ベクトル演算写像、順序写像、汎関数演算写像、固定点演算写像の 1 つ以上を含むことを特徴とする請求項 1 ないし 9 いずれか 1 項に記載のパターン解析方法。

10

【請求項 11】

コンピュータプログラムを含む記憶手段と、該コンピュータプログラムを実行する時に、

少なくとも 1 つの第 1 データを受信する第 1 工程と、

該第 1 データを、データ、写像、およびパターンの記録可能な前記記憶手段上のデータ構造に記録する第 2 工程と、

該第 1 データの構造と解釈に応じて少なくとも 1 つの原始写像を決定する第 3 工程と、

該原始写像を前記データ構造に記録する第 4 工程と、

前記データ構造に記録された少なくとも 1 つの第 2 データを選択する第 5 工程と、

20

該第 2 データから少なくとも 1 つの第 3 データを誘導する複数の手続きから少なくとも 1 つの手続きを選択する第 6 工程と、

該手続きに従い前記第 2 データから前記第 3 データを誘導する第 7 工程と、

前記データ構造に該第 3 データを記録する第 8 工程と

を実行するように配置された処理手段を有し、前記複数の手続きに、前記データ構造に記録された少なくとも 1 つの第 1 写像を所定の方法により選択し該第 1 写像を前記第 2 データに適用することにより前記第 3 データを誘導することを特徴とする第 1 手続きと、

前記第 2 データ内に表現された複数の集合の直積を少なくとも 1 つとり該直積を前記第 3 データ内に表現することにより該第 3 データを誘導することを特徴とする第 2 手続きと

30

を含むことを特徴とするパターン解析システムであって、

前記処理手段が前記コンピュータプログラムを実行する時に、

前記データ構造に記録された少なくとも 1 つの第 4 データ内に所定の方法により少なくとも 1 つの第 1 パターンを探す第 9 工程

をさらに実行するように配置されたパターン解析システム。

【請求項 12】

前記処理手段が前記コンピュータプログラムを実行する時に、

前記第 1 パターンが見つければ該第 1 パターンを前記データ構造に記録する第 10 工程と、

前記データ構造に記録された少なくとも 1 つの第 2 パターンをパターン解析結果として提供する第 11 工程

40

をさらに実行するように配置された請求項 11 記載のパターン解析システム。

【請求項 13】

前記処理手段が前記コンピュータプログラムを実行する時に、

前記第 1 パターンが見つければ該第 1 パターンに対応する理想化されたデータである少なくとも 1 つの第 5 データを所定の方法により生成し前記データ構造に記録する第 12 工程をさらに実行するように配置された請求項 11 ないし 12 いずれか 1 項に記載のパターン解析システム。

【請求項 14】

前記第 12 工程は

前記第 4 データ内に表現された少なくとも 1 つの第 1 確率測度を選択しエントロピーのよ

50

り低い第2確率測度を該第1確率測度から生成し該第2確率測度を前記第5データ内に表現するか、

前記第4データ内に表現された少なくとも1つの第3確率測度を選択し該第3確率測度を集中させて少なくとも1つの第4確率測度を生成し該第4確率測度を前記第5データ内に表現するか、

前記第4データ内に表現された少なくとも1つの第5確率測度を選択し該第5確率測度内の少なくとも1つの確率の集中に各々対応した複数の確率測度を生成し該複数の確率測度を前記第5データ内に表現するか、

前記第4データ内の近似的に繰り返すパターンをより正確に前記第5データ内に繰り返させるかの少なくとも1つにより前記第5データを生成することを特徴とした請求項13記載のパターン解析システム。

10

【請求項15】

前記処理手段が前記コンピュータープログラムを実行する時に、

前記第1パターンが見つければ前記第4データに結び付けられた手続き及び写像に従ってパターン写像を決定し前記データ構造に記録する第13工程をさらに実行するように配置され、前記第8工程は選択された前記手続きと該手続きで写像が使用されていれば該写像を前記第3データと結び付けて前記データ構造に記録することを特徴とした請求項11ないし14いずれか1項に記載のパターン解析システム。

【請求項16】

前記処理手段が前記コンピュータープログラムを実行する時に、

一連の工程を所定の停止条件が満たされるまで繰り返す第14工程をさらに実行するように配置され、該一連の工程は前記第5ないし14工程のうち前記第14工程以外の少なくとも1つを含むことを特徴とした請求項11ないし15いずれか1項に記載のパターン解析システム。

20

【請求項17】

前記複数の手続きに、

前記データ構造に記録された少なくとも1つの第2写像を選択し前記第2データ内に表現された少なくとも1つの第1集合の該第2写像による逆像をとり前記第3データ内に該逆像を表現することにより該第3データを誘導することを特徴とする第3手続きをさらに含むことを特徴とした請求項11ないし16いずれか1項に記載のパターン解析システム。

30

【請求項18】

前記複数の手続きに、

前記第2データ内に表現された少なくとも1つの第2集合の少なくとも1つの部分集合をとり前記第3データ内に該部分集合を表現することにより該第3データを誘導することを特徴とする第4手続きをさらに含むことを特徴とした請求項11ないし17いずれか1項に記載のパターン解析システム。

【請求項19】

前記原始写像が、

恒等写像、定数写像、等号写像、積写像、複数の写像の積写像を与える写像、引戻し演算写像、射影写像、対角写像、置換写像、写像合成写像、評価写像、複数の低位写像を組み合わせて高位写像を与える写像、CURRY写像、論理演算写像、ベクトル演算写像、順序写像、汎関数演算写像、固定点演算写像の1つ以上を含むことを特徴とする請求項11ないし18いずれか1項に記載のパターン解析システム。

40

【請求項20】

記憶手段を備えた処理手段によって実行されたとき、パターン解析をするように配置されたソフトウェア・プログラムが記録されたソフトウェア記録媒体であって、該ソフトウェア・プログラムは

実行されたとき、少なくとも1つの第1データを受信する第1モジュールと、

実行されたとき、該第1データを、データ、写像、およびパターンの記録可能な前記記憶手段上のデータ構造に記録する第2モジュールと、

50

実行されたとき、該第 1 データの構造と解釈に応じて少なくとも 1 つの原始写像を決定する第 3 モジュールと、

実行されたとき、該原始写像を前記データ構造に記録する第 4 モジュールと、

実行されたとき、前記データ構造に記録された少なくとも 1 つの第 2 データを選択する第 5 モジュールと、

実行されたとき、該第 2 データから少なくとも 1 つの第 3 データを誘導する複数の手続きから少なくとも 1 つの手続きを選択する第 6 モジュールと、

実行されたとき、前記手続きに従い前記第 2 データから前記第 3 データを誘導する第 7 モジュールと、

実行されたとき、前記データ構造に該第 3 データを記録する第 8 モジュールと

10

を有し、前記複数の手続きに、

前記データ構造に記録された少なくとも 1 つの第 1 写像を所定の方法により選択し該第 1 写像を前記第 2 データに適用することにより前記第 3 データを誘導することを特徴とする第 1 手続きと、

前記第 2 データ内に表現された複数の集合の直積を少なくとも 1 つとり該直積を前記第 3 データ内に表現することにより該第 3 データを誘導することを特徴とする第 2 手続きと

を含むことを特徴とするソフトウェア記録媒体

であって、

前記ソフトウェア・プログラムが実行されたとき、前記データ構造に記録された少なくとも 1 つの第 4 データ内に所定の方法により少なくとも 1 つの第 1 パターンを探す第 9 モジュール

20

をさらに有することを特徴としたソフトウェア記録媒体。

【請求項 2 1】

前記ソフトウェア・プログラムが

実行されたとき、前記第 1 パターンが見つければ該第 1 パターンを前記データ構造に記録する第 10 モジュールと、

実行されたとき、前記データ構造に記録された少なくとも 1 つの第 2 パターンをパターン解析結果として提供する第 11 モジュールと

をさらに有することを特徴とした請求項 2 0 記載のソフトウェア記録媒体。

【請求項 2 2】

30

前記ソフトウェア・プログラムが

実行されたとき、前記第 1 パターンが見つければ該第 1 パターンに対応する理想化されたデータである少なくとも 1 つの第 5 データを所定の方法により生成し前記データ構造に記録する第 12 モジュール

をさらに有することを特徴とした請求項 2 0 ないし 2 1 いずれか 1 項に記載のソフトウェア記録媒体。

【請求項 2 3】

前記第 12 モジュールは、実行されたとき、

前記第 4 データ内に表現された少なくとも 1 つの第 1 確率測度を選択しエントロピーのより低い第 2 確率測度を該第 1 確率測度から生成し該第 2 確率測度を前記第 5 データ内に表現するか、

40

前記第 4 データ内に表現された少なくとも 1 つの第 3 確率測度を選択し該第 3 確率測度を集中させて少なくとも 1 つの第 4 確率測度を生成し該第 4 確率測度を前記第 5 データ内に表現するか、

前記第 4 データ内に表現された少なくとも 1 つの第 5 確率測度を選択し該第 5 確率測度内の少なくとも 1 つの確率の集中に各々対応した複数の確率測度を生成し該複数の確率測度を前記第 5 データ内に表現するか、

前記第 4 データ内の近似的に繰り返すパターンをより正確に前記第 5 データ内に繰り返させるかの少なくとも 1 つにより前記第 5 データを生成することを特徴とした請求項 2 2 記載のソフトウェア記録媒体。

50

【請求項 2 4】

前記ソフトウェア・プログラムが実行されたとき、前記第 1 パターンが見つければ前記第 4 データに結び付けられた手続き及び写像に従ってパターン写像を決定し前記データ構造に記録する第 1 3 モジュールをさらに有し、

前記第 8 モジュールは選択された前記手続きと該手続きで写像が使用されていれば該写像を前記第 3 データと結び付けて前記データ構造に記録することを特徴とした請求項 2 0 ないし 2 3 いずれか 1 項に記載のソフトウェア記録媒体。

【請求項 2 5】

前記ソフトウェア・プログラムが

実行されたとき、一連のモジュールを所定の停止条件が満たされるまで繰り返し実行する第 1 4 モジュールをさらに有し、該一連のモジュールは前記第 5 ないし 1 4 モジュールのうち前記第 1 4 モジュール以外の少なくとも 1 つを含むことを特徴とした請求項 2 0 ないし 2 4 いずれか 1 項に記載のソフトウェア記録媒体。

10

【請求項 2 6】

前記複数の手続きに、

前記データ構造に記録された少なくとも 1 つの第 2 写像を選択し前記第 2 データ内に表現された少なくとも 1 つの第 1 集合の該第 2 写像による逆像をとり前記第 3 データ内に該逆像を表現することにより該第 3 データを誘導することを特徴とする第 3 手続きをさらに含むことを特徴とした請求項 2 0 ないし 2 5 いずれか 1 項に記載のソフトウェア記録媒体。

20

【請求項 2 7】

前記複数の手続きに、

前記第 2 データ内に表現された少なくとも 1 つの第 2 集合の少なくとも 1 つの部分集合をとり前記第 3 データ内に該部分集合を表現することにより該第 3 データを誘導することを特徴とする第 4 手続きをさらに含むことを特徴とした請求項 2 0 ないし 2 6 いずれか 1 項に記載のソフトウェア記録媒体。

20

【請求項 2 8】

前記原始写像が、

恒等写像、定数写像、等号写像、積写像、複数の写像の積写像を与える写像、引戻し演算写像、射影写像、対角写像、置換写像、写像合成写像、評価写像、複数の低位写像を組み合わせ高位写像を与える写像、CURRY 写像、論理演算写像、ベクトル演算写像、順序写像、汎関数演算写像、固定点演算写像の 1 つ以上を含むことを特徴とする請求項 2 0 ないし 2 7 いずれか 1 項に記載のソフトウェア記録媒体。

30

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明はデータ解析、特にパターンが発見できるようにデータを配置する方法と装置に関する。

【背景技術】

【0 0 0 2】

データ管理、データ処理、そしてデータ解析は現代生活及び仕事の上での偏在的要素となっている。科学的、医学的、工学的、そして商業的用途のための膨大なデータの流れの開発、管理、そして保管は、巨大産業となった。バイオテクノロジー、金融、画像、その他のデータのソース、及び需要は、急速に増大している。必ずしもどれが興味ある現象に関係あるか知らないまま、体系的に多数の測定が行われ、大量のデータが自動的に集められている。

40

【0 0 0 3】

したがって、適切な情報を巨大なデータの山から捻り出し、干し草の山の中の針を探すことが益々重要になっている。これは、現在データ解析で使われている多くの技術の背後にある古い仮定と重大な違いを持つ。これらの技術の多くは、例えば科学的知識により前もって丁度適切な変数を測定するなどして精選された、少数の変数を対象にすることを前提

50

にしている。

【 0 0 0 4 】

これらの技術で使われている基本的な方法論は、常に適用可能では既になくなっていく。データ解析における以前の方法の基礎をなす理論はデータ要素の数が個別データの次元より遥かに大きいことを前提としている。しかし、今日ではしばしば、データの次元はデータ要素の数より遥かに大きい。そのような場合はもう例外ではなくある意味で常態である。多くの種類の事象について、その事象を数量化する潜在的に非常に多数の測定可能な項目が存在し、その事象の例の数は比較的少ない。一例を挙げれば、多数の遺伝子と、ある遺伝病に罹患する比較的少数の患者の例がある。他例には画像がある。画像は軽く百万次元（画素）を持ちうるが、一組の解析すべきデータとして百万個の画像を処理することは稀である。

10

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 0 5 】

従って、高次元データがより効果的に解析されうるように与えられたデータを配置する方法と装置を与えるのが本発明の目的である。与えられたデータ内によりうまくパターンを見つけられるようにそのデータを配置する方法を与えるのも本発明の目的である。

【 課題を解決するための手段 】

【 0 0 0 6 】

本方法は与えられたデータ内にパターンが発見できるようにそのデータを配置することを可能とする。データを特徴づける写像及びそれが所属する集合を利用して、本方法は多数の「データ項目」を比較的少数の入力データ項目から作り、統計的その他の既存のデータ解析手法を適用することを可能にする。データ又はその一部から写像の集合が決定される。次に、既存の写像を組合せるか、ある種の変換を写像に加えることにより、新しい写像が生成される。次に、それらの写像をデータに適用した結果が調べられパターンが探される。例えば、本発明のある実施形態においては、特定の結果データあるいはデータの集合の頻度が調べられる。場合により随意に、ある強いパターンが選ばれ、理想化され、そのパターンを反映したデータを見つけるために伝播され戻される。すなわち、本発明のパターン解析方法は記憶手段と処理手段を備えた情報処理システムにより実行されるパターン解析方法であって、少なくとも1つの第1データを受信する第1工程と、該第1データを前記記憶手段に記録する第2工程と、該第1データの構造と解釈に応じて少なくとも1つの原始写像を決定する第3工程と、該原始写像を前記記憶手段に記録する第4工程と、前記記憶手段に記録された少なくとも1つの第2データを選択する第5工程と、該第2データから少なくとも1つの第3データを誘導する複数の手続きから少なくとも1つの手続きを選択する第6工程と、該手続きに従い前記第2データから前記第3データを誘導する第7工程と、前記記憶手段に該第3データを記録する第8工程とを有し、前記複数の手続きに、前記記憶手段に記録された少なくとも1つの第1写像を所定の方法により選択し該第1写像を前記第2データに適用することにより前記第3データを誘導することを特徴とする第1手続きと、前記第2データ内に表現された複数の集合の直積を少なくとも1つとり該直積を前記第3データ内に表現することにより該第3データを誘導することを特徴とする第2手続きとを含むことを特徴とする。また、前記第9工程は前記第4データのエンタロピーに従って前記第1パターンを探すこととすることもできる。また、前記第1パターンが見つければ該第1パターンを前記記憶手段に記録する第10工程と、前記記憶手段に記録された少なくとも1つの第2パターンをパターン解析結果として提供する第11工程をさらに有することもできる。また、前記第1パターンが見つければ該第1パターンに対応する少なくとも1つの第5データを所定の方法により生成し前記記憶手段に記録する第12工程をさらに有することもできる。また、前記第12工程は、前記第4データ内に表現された少なくとも1つの第1確率測度を選択しエンタロピーのより低い第2確率測度を該第1確率測度から生成し該第2確率測度を前記第5データ内に表現するか、前記第4データ内に表現された少なくとも1つの第3確率測度を選択し該第3確率測度を集中させて

20

30

40

50

少なくとも1つの第4確率測度を生成し該第4確率測度を前記第5データ内に表現するか、前記第4データ内に表現された少なくとも1つの第5確率測度を選択し該第5確率測度内の少なくとも1つの確率の集中に各々対応した複数の確率測度を生成し該複数の確率測度を前記第5データ内に表現するか、前記第4データ内の近似的に繰り返すパターンをより正確に前記第5データ内に繰り返させるかの少なくとも1つにより前記第5データを生成することとすることもできる。また、前記第1パターンが見つければ前記第4データに結び付けられた手続き及び写像に従ってパターン写像を決定し前記記憶手段に記録する第13工程をさらに有し、前記第8工程は選択された前記手続きと該手続きで写像が使用されていれば該写像を前記第3データと結び付けて前記記憶手段に記録することとすることもできる。また、一連の工程を所定の条件が満たされるまで繰り返す第14工程をさらに有し、該一連の工程は前記第5ないし14工程のうち前記第14工程以外の少なくとも1つを含むこととすることもできる。また、以上において、前記複数の手続きに、前記記憶手段に記録された少なくとも1つの第2写像を選択し前記第2データ内に表現された少なくとも1つの第1集合の該第2写像による逆像をとり前記第3データ内に該逆像を表現することにより該第3データを誘導することを特徴とする第3手続きをさらに含むこともできる。また、前記複数の手続きに、前記第2データ内に表現された少なくとも1つの第2集合の少なくとも1つの部分集合をとり前記第3データ内に該部分集合を表現することにより該第3データを誘導することを特徴とする第4手続きをさらに含むこともできる。また、前記原始写像が、恒等写像、定数写像、等号写像、積写像、複数の写像の積写像を与える写像、引戻し演算写像、射影写像、対角写像、置換写像、写像合成写像、評価写像、複数の低位写像を組み合わせて高位写像を与える写像、CURRY写像、論理演算写像、ベクトル演算写像、順序写像、汎関数演算写像、固定点演算写像の1つ以上を含むこととすることもできる。

10

20

【0007】

本発明のいくつかの側面の基本的理解を与えるために、以下に本発明の単純化された要約を示す。この要約は本発明の広範囲にわたる概観ではなく、本発明の鍵となるあるいは決定的な要素を指示することや、本発明の範囲を定めることも意図していない。その唯一の目的は、後のより詳細な記述への前触れとして、本発明のいくつかの概念を単純化した態様で示すことにある。

データ

30

【0008】

図1にデータ内にパターンを発見する方法のフローチャートを示す。本方法によれば、まず解析されるべきデータが受信される(101)。最も一般的なデータの形態は、遍在する情報処理システムや機器で使われるようなビットの列である。データは通常、何らかの構造と解釈を持つ。例えば、データのある部分は8ビットのグループ毎に一文字と解釈されるテキストデータかもしれない。他の部分は32ビット整数や64ビット浮動小数点数を表すかもしれない。あるいは単一のビットが「yes」又は「no」と解釈されるかもしれない。遺伝子配列を表すデータでは、2ビットでヌクレオチド中の塩基(A, G, C, Tのどれか)を表すかもしれない。データはそれぞれが一組の情報を表すいくつかのレコードに分割されている場合もある。例えば画像データは画素数(幅と高さ)を表す2つの整数と、各画素の色を表す整数の列からなるかもしれない。

40

【0009】

表記

【0010】

以下では、データをもう少し抽象的に取り扱う。整数は、それを表すのに何ビット使われているかが整数と呼ぶ。同様に、浮動小数点数のことは実数と呼び、「yes」と「no」のように二者択一を表すデータは全てブール値と呼ぶ。さらに一般に、以下では各種の集合と写像について言及する。

【0011】

集合は要素の集まりである。例えば、整数の集合Zは全ての整数を要素に持つ集合である

50

。ブール値の集合 `bool` は `true` と `false` の2つの要素しか持たない。集合はその全ての要素を「`{ }`」内に列挙して表記されることもある。例えば `bool = { true, false }` のように。表記 `a ∈ A` は `a` が集合 `A` の要素であることを表す。もし集合 `B` の全ての要素がもう1つの集合 `A` の要素でもあるならば、`B` は `A` の部分集合であり、これを `A ⊇ B` (又は `B ⊆ A`) と表記する。二つの集合 `A` と `B` は、もし `A ⊆ B` かつ `B ⊆ A` ならば等しい (`A = B` と表記)。 `A` の部分集合 `B` は、もし `A ⊇ B` ならば真部分集合である。

【0012】

これらの表記を使うことは、本発明が実際に集合という数学的概念を扱うことを意味しない。それは本方法を簡潔に、(これらの表記が概念の記述に、しばしばあまり厳密でなく、使われる) 関連技術分野で熟練した者によく知られた表記で記述するためである。例えば、`Z` のように無限個の要素を持つ集合があり、また(実数のように) 正確に指定するには無限の精度を要する要素を持つ集合もあるが、それらは有限の存在である情報システムで日常的に扱われている。これは、普通はそのような集合の有限個の要素しか、各仕事に必要なからである。また集合は時に記号的に処理され、または近似されることもある。集合や写像を表現し操作するこれらや他の手法は関連技術分野である計算機科学では良く知られている。`SETL` や `MIRANDA` 等のある種のプログラム言語は集合を言語プリミティブとして持つほどである。また、ここで使われる集合や写像の概念は `ML` や `HASKELL` のような型付き関数型言語における型と写像の概念に非常に近い。従って、関連技術分野における通常的能力を持つ者は適切な手法を使ってここに開示される本方法を実現することができるであろう。

【0013】

集合 `A` と `B` に対して、「`A → B`」は `A` から `B` への写像の集合を表す。写像とは与えられた集合の各要素に一意の対象を結びつける方法である。つまり `A` から `B` への写像とは、`A` の各要素 `a` にただ一つの `B` 内の対象 `f(a)` を与えるような関数のことである。そのような状態は時に「`f` は `a` を `f(a)` に送る(あるいは写像する)」と記述される。表記「`f : A → B`」は `f` が集合 `A` から集合 `B` への写像であること、即ち `f` が `A → B` の要素であることを意味する。写像 `f : A → B` に対して、`A` は `f` の定義域と呼ばれる。

【0014】

集合 `A` について、`idA : A → A` は `A` の各要素 `a` をそれ自身に送る恒等写像を表す。

【0015】

集合 `A` と `B` について、定数写像 `const : A → B` は `const(a) = b` で定義される。つまり、`A` の要素 `a` について、`const(a) : B` は `B` のどんな要素 `b` も `a` に送る写像である。

【0016】

`B` が `A` の部分集合である時、包含写像 `incl : B → A` は `incl(b) = b` で定義される。

【0017】

2つの集合 `A` と `B` について、`A × B` はこの2集合の直積、つまり順序対 `(a, b)` (`a` は `A` に、`b` は `B` に属する) の集合を表す。同様に、`A × B × C` は3集合 `A`、`B`、`C` の直積を表し、以下同様である。一般に別の集合 `I` でインデックスされた任意の集合族 `Ai` の直積は、 $\prod_{i \in I} A_i$ または、`Ai` が全て等しい時には `AI` で表される。 $\prod_{i \in I} A_i$ の要素は `(ai)i ∈ I` で表される。ここで各 `ai` は `Ai` の要素である。有限個の要素を持つ標準集合を次のように表記する。`Z1 = { 1 }`、`Z2 = { 1, 2 }`、`...`、`Zn = { 1, ..., n }`。以下では、`A × B` は `I = Z2`、`A1 = A`、`A2 = B` のときの $\prod_{i \in I} A_i$ の略記と理解されたい。同様に、`A × B × C` は `I = Z3`、`A1 = A`、`A2 = B`、`A3 = C` のときの $\prod_{i \in I} A_i$ の略記であり、以下同様である。

【0018】

写像 `f : A → B` は、各 `a ∈ A` について `f` の `a` 番目の要素を `f(a)` と考えることによって、`BA` 即ち `A` でインデックスづけされた `B` のコピーの直積の要素と考えられる。したがって、`A → B` はここでは `BA` の別名とみなされる。

【0019】

特別な集合 $unit$ が定義される。それはただ1つの要素を持つ。 $unit$ によって、集合 A の任意の要素 a を、 $unit$ の唯一の要素を a に送る写像 $a : unit \rightarrow A$ とみなすことができる。写像に対してのみに適用可能な写像または操作を集合 A の(写像でない)普通の要素に適用するために、本発明が自動的にこの変換を実行することもある。 A^{unit} あるいは $unit \rightarrow A$ という形の集合は A と同一視される。

【0020】

写像 $f : A \rightarrow B$ と B の要素 b について、 f による b の逆像 $f^{-1}(b)$ とは、 f により b に送られる A の要素からなる A の部分集合である。 B の部分集合 C の f による逆像 $f^{-1}(C)$ とは、 f により C の要素に送られる A の要素からなる A の部分集合である。

10

【0021】

ある種の写像は再帰的に定義される。つまり、再帰的に定義された写像はその定義にそれ自身を使用する。例えば、階乗関数 $fac : \mathbb{N} \rightarrow \mathbb{N}$ は自然数 n を、もし n が1ならば1に、それ以外ならば $fac(n)$ の n 倍に送る関数として定義される(ここで \mathbb{N} は自然数の集合 $\{1, 2, 3, \dots\}$ を表す)。

【0022】

引戻し

【0023】

2つの積集合 $\prod_{i \in I} A_i$ と $\prod_{j \in J} B_j$ について、全ての $j \in J$ について $A_{h(j)} = B_j$ である写像 $h : J \rightarrow I$ が存在する時、対応する引戻し $h^* : \prod_{j \in J} B_j \rightarrow \prod_{i \in I} A_i$ が $(h^*(\langle a_i \rangle_{i \in I}))_j = a_{h(j)}$ で定義される。この写像の特別な場合に以下がある。

20

【0024】

[PB1] I の任意の部分集合 J について、 $h = incl : J \rightarrow I$ とした $h^* : \prod_{j \in J} A_j \rightarrow \prod_{i \in I} A_i$ は射影写像を定義する。例えば直積 $A \times B$ について、自然な射影がある：

- $proj_A : A \times B \rightarrow A$ [$proj_A(a, b) = a$]
- $proj_B : A \times B \rightarrow B$ [$proj_B(a, b) = b$]

写像 $proj_A$ は $h^* : \prod_{j \in Z_2} A_j \rightarrow \prod_{i \in Z_2} A_i$ で $A_1 = A$ 、 $A_2 = B$ 、 $h = incl : Z_1 \rightarrow Z_2$ としたものと同一である。

30

【0025】

[PB2] 同じ集合 n 個のコピーの直積 $A \times A \times \dots \times A$ について、対角写像 $diag : A \times A \times \dots \times A \rightarrow A \times A \times \dots \times A$ が $diag(a) = (a, a, \dots, a)$ で定義される。これは $h^* : \prod_{j \in Z_n} B_j \rightarrow \prod_{i \in Z_n} A_i$ で $A_1 = A$ 、 $B_j = A$ とし、 $h : Z_n \rightarrow Z_1$ を $Z_n = \{1, \dots, n\}$ の全ての j について $h(j) = 1$ で定義したものと同一である。

【0026】

[PB3] 直積 $A \times B$ について、 (a, b) を (b, a) に送る交換写像 $A \times B \rightarrow B \times A$ がある。同様に任意の数の集合の直積について、成分の順序を変える置換写像がある。これは $h^* : \prod_{j \in Z_n} B_j \rightarrow \prod_{i \in Z_n} A_i$ で h を置換写像とし $Z_n = \{1, \dots, n\}$ の全ての j について $B_j = A$ としたものと同一である。

40

【0027】

[PB4] 2つの写像 $f : A \rightarrow B$ と $g : B \rightarrow C$ について、合成写像 $g \circ f : A \rightarrow C$ が A 内の a について $g \circ f(a) = g(f(a))$ で定義される。これも引戻しの特別な場合である。これをみるには全ての C_b と C_a を C と等しくして $g : C^B \rightarrow C^A$ 、 $f : C^A \rightarrow C^B$ であり $g \circ f : C^A \rightarrow C^A$ であることを思い出されたい。

【0028】

[PB5] 集合 A と B 、 A 内の a について $const(a) : B \rightarrow A$ は B 内の任意の b を a に送る写像である。 $J = Z_n$ とした定数写像 $const(a) : J \rightarrow A$ とその引戻し $const(a)^* : \prod_{j \in J} B_j \rightarrow \prod_{i \in A} B_i$ を考える。それは写像 $f : A \rightarrow B$ を、その a

50

での値 $f(a)$ B に写像する。これは、 $ev(f, a) = f(a)$ で定義される、写像の値を評価する写像 $ev : (A \rightarrow B) \times A \rightarrow B$ を定義する。

【0029】

統計

【0030】

本発明においては、データを確率測度（確率分布）のような統計として表現すること、あるいはもっと一般に、データの相対頻度を処理することが、特に有用である。一般に、集合 A について、 A 上の確率測度 Pr は A の（事象と呼ばれる）部分集合 B に対して 0 と 1 の間の実数 $Pr(B)$ を与える。データを確率測度で表すとは以下を意味する。もしあるデータが集合 A の単一の要素 a であるなら、それは A の事象 B が a を含むときには $Pr(B) = 1$ を与え、それ以外するとき $Pr(B) = 0$ を与える確率測度として表現されうる。あるいはそれは a を中心としたガウス分布のような、概算測度としても表現されうる。同じ集合に属する多くのデータ点があるときには、 A に含まれる全てのデータ点に対する B に含まれるデータ点の比を与える、単純な係数測度 $Pr(B)$ として表現されるかもしれないし、あるいは再び、ガウス混合分布や *Parzen Window* の手法のような概算測度としても表現されうる。情報システムにおけるそのような確率測度の処理及びシミュレーションのための種々の手法が、関連技術分野ではよく知られている。後述のある実施形態においては、頻度係数と呼ばれる具体的な方法が使われる。このように確率測度を使うとき、各週強情の標準測度が必要に応じて使われる。これは、一様分布のように、その集合の、特徴のないデフォルト状態を表す確率測度である。

原始写像

【0031】

次に、そのデータかそのデータの一部からの写像の集合が決定される（102）。これらの写像は原始写像と呼ばれる。原始写像に含まれる写像は集合上に定義される標準写像の一つかもしれない。例えば、整数の集合 Z には、ある整数をその次に数に送る、自身への写像がある。集合 Z にはまた加法もある。それは $Z \times Z$ から Z への写像として表現されるのだが、これも原始写像の集合に加えられるかもしれない。このように加法写像は $Z \times Z$ 内の (i, j) を Z 内の $i + j$ に送る。従って、データの一部が一つあるいは複数の整数を表していれば、その整数の次の数を与える写像あるいはそれらの整数の和を与える写像が原始写像に含められるかもしれない。ある種の集合にはそれらの間に自然な写像を持つ。例えば、任意の集合 A について、等しさという概念は $A \times A$ からブール値の集合 $bool = \{true, false\}$ への写像を定義する。つまり、 $A \times A$ 内の (u, v) に対してその写像は $u = v$ であるときに限り $true$ を与える。同様に、ある種の集合には順序の概念があり、写像と考えることができる。例えば整数の集合 Z に、 $Z \times Z$ 内の (i, j) に $i < j$ のときに限り $true$ を与える $Z \times Z$ から $bool$ への順序写像がある。

【0032】

以下に集合に自然に随伴し、原始写像の集合に含められるかもしれない写像のいくつかを列挙する。ここで R は実数の集合を表す。

【0033】

[PM I] 任意の集合 A は次の原始写像を持つ：

・恒等写像： $id_A : A \rightarrow A$ [$id_A(a) = a$]

・定数写像： $const : A \rightarrow (B \rightarrow A)$ [$const(a)(b) = a$]（任意の集合 B について）

【0034】

[PM II] 等しいかどうか簡単に決定できる集合 A について、等号写像：

・ $eq_A : A \times A \rightarrow bool$ [$a = b$ なら $eq_A(a, b) = true$ 、それ以外 $false$]

【0035】

[PM III] 2つの写像 $f : A \rightarrow B$ と $g : C \rightarrow D$ について、積写像 $f \times g : A \times C \rightarrow B \times D$ が $f \times g((a, c)) = (f(a), g(c))$ で定義される。これは原始写像

10

20

30

40

50

$mp : (A \ B) \times (C \ D) \ (A \times C \ B \times D)$ を定義する。

【0036】

[PM IV] 写像に対する引戻し: $pullback : (J \ D \ (\ i \ I \ A \ i \ j \ J \ B \ i))$ 。これは写像を別の写像に送る。これの特別な場合には射影写像 [PB1]、対角写像 [PB2]、置換写像 [PB3]、写像合成写像 [PB4]、評価写像 [PB5] が含まれる。

【0037】

[PM V] 低位写像の組み合わせ。K をインデックスの集合とし、各 $k \in K$ について I_k もインデックスの集合とする。 $k \in K$ について既知の写像 $f_k : \prod_{i \in I_k} A_{k,i} \rightarrow \prod_{i \in I_k} B_{k,i}$ があり、もう一つのインデックス集合 J と、 $A_{k,i} = A_{m,j}$ のとき $h_k(i) = h_m(j)$ であるような写像 $h : \prod_{j \in J} A_{k,i} \rightarrow \prod_{j \in J} B_{k,i}$ と $h : \prod_{j \in J} A_{k,i} \rightarrow \prod_{j \in J} B_{k,i}$ を、F は K 内の全ての k についての f_k の積集合として、 $L = \bigcup_{k \in K} I_k$ はインデックス集合 I_k の共通部分のない和集合として、そして h は I_k 上で h_k と一致するように、それぞれ定義する。すると、h の引戻し $h^* : \prod_{j \in J} A_{k,i} \rightarrow \prod_{j \in J} B_{k,i}$ と F を合成すると新しい写像 $F \circ h^* : \prod_{j \in J} A_{k,i} \rightarrow \prod_{j \in J} B_{k,i}$ が定義される。

10

【0038】

[PM VI] $curry$ 写像 $curry : (A \times B \rightarrow C) \rightarrow (A \rightarrow (B \rightarrow C))$ は、写像 $f : A \times B \rightarrow C$ を、写像 $curry(f) : A \rightarrow (B \rightarrow C)$ に送るが、これは A 内の a を $curry(f)(a)(b) = f(a, b)$ で定義される写像 $curry(f)(a) : B \rightarrow C$ に送る。逆の操作は $uncurry$ 写像 $uncurry : (A \rightarrow (B \rightarrow C)) \rightarrow (A \times B \rightarrow C)$ で、これは写像 $g : A \rightarrow (B \rightarrow C)$ を、 $(a, b) \in A \times B$ を $g(a)(b)$ に送る別の写像 $uncurry(g) : A \times B \rightarrow C$ に送る。これは計算機科学ではよく知られている。

20

【0039】

[PM VII] 各種の論理演算がある: NOT: $bool \rightarrow bool$ 、AND: $bool \times bool \rightarrow bool$ 、OR: $bool \times bool \rightarrow bool$ 等。

【0040】

[PM VIII] R を含む任意のベクトル空間は次の自然な写像を持つ:

- ・ (加法) $Add_V : V \times V \rightarrow V$ [$Add_V(u, v) = u + v$]
- ・ (実数との積) $Mult_V : R \times V \rightarrow V$ [$Mult_V(a, v) = a \cdot v$]
- ・ (減法) $Sub_V : V \times V \rightarrow V$ [$Sub_V(u, v) = u - v$] (これは加法と -1 倍によって定義できるが、後の記法の簡略化のためにここに含める。)
- ・ (長さ) $Len_V : V \rightarrow R$ [$Len_V(v) = \text{ベクトル } v \text{ の長さ}$]
- ・ 別のベクトル空間でパラメータづけられた種々の線形変換: $LT : V \times U \rightarrow W$
- ・ 別のベクトル空間でパラメータづけられた種々の双線形、trilinear、... 等々の形式:

30

・ $LF : V \times U \rightarrow R$

・ $BF : V \times V \times U \rightarrow R$

・ $TF : V \times V \times V \times U \rightarrow R$

40

【0041】

[PM IX] R は順序の概念を持つ:

- ・ $Ord_R : R \times R \rightarrow bool$ [$a < b$ なら $Ord_R(a, b) = true$ 、それ以外は $false$]

【0042】

[PM X] ユークリッド空間 E は 2 点間のベクトルの概念を持つ:

- ・ $Diff_E : E \times E \rightarrow V$ (V は同次元のベクトル空間)

【0043】

[PM XI] R の部分集合 A 上の実数値関数のある種の集合 U (つまり U は A \rightarrow R の部分集合) について、微分写像 $Der : U \rightarrow (A \rightarrow R)$ は関数とその導関数 (微分) に送る

50

。実ベクトル空間の間の写像の様々な微分をとる同様な写像がある。さらに一般に、原始写像として加えられるかもしれないよく知られた数学的変換は他にもある（例えばフーリエ変換）。

【0044】

[PM XII] 固定点演算。写像 $f: A \rightarrow A$ について、固定点演算子 $\text{Fix}: (A \rightarrow A) \rightarrow A$ はその写像のある固定点を与える。つまり、 $a = \text{Fix}(f)$ は $f(a) = a$ であるような A の要素である。これは、再帰的に定義される写像を定義するのに使える。例えば、上述の階乗写像 $\text{fac}: \mathbb{N} \rightarrow \mathbb{N}$ を再帰的でない写像から得ることができる。写像 $f: \mathbb{N} \rightarrow \mathbb{N}$ を別の写像 $F(f): \mathbb{N} \rightarrow \mathbb{N}$ に送る写像 $F: (\mathbb{N} \rightarrow \mathbb{N}) \rightarrow (\mathbb{N} \rightarrow \mathbb{N})$ を次のように定義する。 $F(f)$ は自然数 n を $n = 1$ なら 1 に、それ以外なら $f(n - 1)$ の n 倍に送る。このとき、 $\text{Fix}(F)$ が階乗写像である。固定点演算は全ての写像に適用可能ではないかもしれないことに注意せよ。

10

【0045】

原始写像はまた、表現されたデータにもっと特有のものであるかもしれない。もしデータ中のある整数がある人の課税所得を表すなら、その所得に対する税額を与える写像も、アプリケーションの必要に応じて、原始写像として含められるかもしれない。

誘導データと写像

【0046】

次のステップ(103)では、そのデータと原始写像をもとに、他のデータや写像が生成される。これらの生成法のうちのいくつかは以下の通り。

20

・ 2つ以上の集合から積集合を作られ得る。積集合上の確率測度は元の集合上のものから誘導され得る。

・ データは写像によって送られ得る。確率測度は写像によって誘導され得る。

・ 集合の写像による逆像がとられ得る。

・ データは部分集合に制限され得る。確率測度も部分集合に制限され得る。

・ 写像を別の写像に送る写像が適用されて、新しい写像が作られ得る、例えば：

・ 2つの写像 $f: A \rightarrow B$ と $g: C \rightarrow D$ から、積写像 $f \times g: A \times C \rightarrow B \times D$ が $f \times g((a, c)) = (f(a), g(c))$ で定義される。([PM III] 参照)

・ 2つの写像 $f: A \rightarrow B$ と $g: B \rightarrow C$ から、写像 $g \circ f: A \rightarrow C$ が A 内の a について $(g \circ f)(a) = g(f(a))$ で定義される。([PM IV] 参照)

30

・ より高位の写像、つまり引数のより多い写像は、多くの対象の間の関係を定義するため重要である。写像を組み合わせるより高位の写像に導くことは、原始写像の殆どは多くても2つの引数しか持たないことから、特に重要である。このように、[PM V]の原始写像は重要である。それは上述した写像の写像を適用する特殊な場合に過ぎないが、ここで例を使って簡単に説明する価値がある。 $f: A \times A \rightarrow B$ を写像とする。高位の写像を作るために、まず積写像を作る： $f \times f: A \times A \times A \times A \rightarrow B \times B$ 。しかしこれは同じことを2回やっているだけだから、あまり多くの新情報をもたらさない。しかし、 $g(a, b, c) = f \times f(a, b, b, c)$ で定義される $g: A \times A \times A \rightarrow B \times B$ は3つの引数の間に新しい関係を定義する。これが、[PM V]の原始写像が適用されたときにこの場合に起こることである。

40

【0047】

上に列挙したように、本方法の様々な段階で新しいデータと写像を生成するための、方法と源の選び方はたくさんある。アプリケーションと、既に見つかったデータと写像を基に有用なパターンを見つける可能性がよりよくなるように、生成されるデータと写像を選ぶための計画があるべきである。一般に、パターン写像(下記参照)とされた写像は、新しい写像の構成要素として使われるより強い傾向を持つべきである。また、なにかのパターンが見つかった集合は源の集合としてより頻繁に使われるべきである。本発明の実施例で使われている一つの方法を後述する。

パターン

【0048】

50

次のステップ(104)では、生成された様々なデータと写像の中にパターンが存在するかどうか調べられる。これは、繰り返されたデータを見つけたり、確率測度の低いエントロピーのような統計的に有意な条件を追求したり、比較的少数の要素への確率の集中を検出するなどの、パターン発見のための従来手法のいずれでもを使ってなされる。以下では、その中にパターンが見つかったデータをパターンデータと呼ぶ。

【0049】

パターンデータは元のデータと生成されたデータに何かの写像を適用した結果である。これらの写像を以下ではパターン写像と呼ぶ。パターン写像はパターン解析に重要である。例えばもし写像をデータに適用した結果がおおまかに繰り返すパターンであるとか、あるいはある確率測度からある写像で誘導された確率測度が低いエントロピーを持っているなら、これらの写像は元のデータを何らかの面で特徴付けている。このパターン写像は類似のデータ中に同じ特徴があるかどうか調べるために適用するのに有用であろう。様々なパターン写像の組合せは、元の集合や各中間段階のデータを特徴付けるかもしれない。

10

【0050】

パターンの存在を決定するとき、写像そのものから来るものを考慮に入れねばならない。つまり、もし写像そのものがパターンを作るなら、そのパターンはデータの特徴を表さない。例えば、上述のエントロピーは、何のパターンも持たない何か(例えばパターン写像の定義集合上の標準確率分布等)に同じパターン写像を適用した結果と相対的に評価しなければならない。

バックトラック

20

【0051】

場合により随意に、次のステップ(105)では、本方法は前ステップで見つかったパターンデータをとってそのパターンに対応する「理想的」データを生成し得る。まず、(パターンデータが見つかったのと)同じ集合内に、パターンデータを修正することで新しいデータが作られるかもしれない。もしそのパターンデータが、生成された集合上のエントロピーの低い確率測度として見つけられたならば、さらに低いエントロピーを持った理想化された確率測度とその集合上に導入されるかもしれない。そして、パターン写像を通してその理想化された確率測度を誘導する確率測度が見つけれられるかもしれない。もし確率の集中が観察されたのなら、理想化はそれをもっと集中するかもしれない。また、もし比較的少数の集中しかないのなら、それぞれ1つの集中を持つ複数の確率測度が、新しいパターンデータとして作られるかもしれない。おおまかに繰り返すパターンは正確に繰り返すパターンにされるかもしれない。

30

【0052】

それから、理想化されたパターンの、対応するパターン写像による逆像がとられるかも知れない。元のデータの入っていた集合まで遡る上での中間段階の集合の中の可能なデータの集合がこうして同定される。これは、そのデータがパターン写像によって理想化されたパターン内に送られたときにtrueを与えるその集合上の述語論理を作ることによって実装され得る。また、元のデータのこの集合内にある部分(つまり、対応する述語論理にtrueを与えられる部分)は特に重要である。なぜならこの部分的データは他の写像によって前へ送られて他に何かパターンが現われるかどうか調べられ得るからである。

40

【0053】

このようにしてパターンを持つ可能なデータの集合が同定できる。十分多くのパターン使い、そのような逆像の共通部分をとることで、可能なデータの小さな集合あるいはただ1つのデータさえ見つかるかもしれない。

【0054】

次のステップ(106)では、望ましいデータが出力される。これは見つかったパターンや、それらに対応する「純粋な」データを含むかもしれない。

【0055】

最後に、プロセスの停止条件が調べられ(107)、もし条件に合わなければプロセスは繰り返す。

50

【図面の簡単な説明】

【0056】

【図1】データ内にパターンを見つける方法のフローチャートを示す。

【図2】探索アルゴリズムのフローチャートを示す。

【図3】データ構造FCと、FC内で使われる部分構造を図式的に表す。

【図4】理想化プロセスのフローチャートを示す。

【発明を実施するための最良の形態】

【0057】

以下の記述では、本発明の完全な理解を与えるために、説明の目的で多数の特定細部が提示される。しかし、関連技術分野で熟練した者には、本発明がそれらの特定細部なしでも実施可能であることが明確であり得る。他の場合には、本発明の記述のために、よく知られた構造や装置がブロックダイアグラム中に示される。本発明は様々な形態のハードウェア、ソフトウェア、ファームウェア、特殊用途プロセッサ、あるいはそれらの組合せによって実装され得ることが理解されるべきである。好ましくは、本発明はプログラム記憶装置に有体的に有形成されたアプリケーションプログラムのソフトウェアとして実装されるべきである。そのアプリケーションプログラムは、任意の適当なアーキテクチャからなる機械に読み込まれ、実行され得る。好ましくは、その機械は、1つあるいは複数の中央処理装置(CPU)、ランダムアクセスメモリ(RAM)、入出力(I/O)インタフェースのようなハードウェアを持つコンピュータプラットフォーム上に実装されるべきである。そのコンピュータプラットフォームはまたオペレーティングシステムとマイクロ命令コードを含む。本明細書に記述された様々なプロセスや関数は、オペレーティングシステムにを通して実行されるそのマイクロ命令コードか、あるいはアプリケーションプログラム、あるいはそれらの組合せであるかもしれない。加えて、追加のデータ記憶装置や印刷装置など、他の様々な周辺装置がそのコンピュータプラットフォームに接続されるかもしれない。さらに理解されるべきことは、付随する図に描かれたシステム構成要素と方法ステップの一部は好ましくはソフトウェアに実装されるべきなので、本発明がプログラムされる態様に依存して、システム構成要素(あるいは方法ステップ)の間の実際の接続は異なるかもしれないことである。本明細書に記述された本発明の教示によれば、関連技術分野の技術者は、これらのあるいは類似した本発明の実装あるいは配置を企図することができるであろう。

データ

【0058】

ここでは、データを分析するための本発明の実施例を提示する。明快さのために、関連技術分野の技術者にはよく知られる一定の抽象性が維持される。例えば、集合や写像は、情報システム上のデータとして表現、又は情報システム上のデータによって近似される。

【0059】

頻度あるいは確率が本発明で以下に操作されるかを描写するため、頻度計数というデータ構造がここに開示される。それは集合上の単純な計数確率測度をモデル化するための具体的方法である。本実施例では、全てのデータはある集合上の頻度計数として表される。

【0060】

以下では、任意の集合Aについて、A上の頻度計数とは、Aの要素とその数を捉えるデータを意味する。それは、Aのいかなる要素も2度以上現れないような $A \times N$ の部分集合として扱われる。ここで $N = \{1, 2, 3, \dots\}$ つまり自然数の集合である。A上の頻度計数の集合は $Freq(A)$ で表される。従ってA上の頻度計数即ち $Freq(A)$ の要素Fは、Aの要素aと自然数nの組 (a, n) の集合Fであって、もし (a, n) を含めば (a, m) の形の他のどんな要素も含まないようなものである。頻度計数内のこれらの組は以下では粒子と呼ばれる。Aの要素aとA上の頻度計数Fについて、 $count_F(a)$ と記述されるaの計数とは、もしF内に (a, n) の形の要素があればnで、なければ0で定義される。 $mass(F)$ すなわちFのマスは、A内の全てのaについての $count_F(a)$ の和として定義される。そしてaの確率 $P_F(a)$ は、 $count_F(a)$

10

20

30

40

50

a) を $mass(F)$ で割ったものと定義される。 F の台 $supp(F)$ は、 $count_F(a) > 0$ である a からなる A の部分集合と定義される。 F のエントロピー $H(F)$ は、 $supp(F)$ 内の全ての a についての和、 $-\sum_{a \in supp(F)} P_F(a) \log_2 P_F(a)$ で定義される。

【0061】

後の参考のために次に注意すべきである。

【0062】

[FC I] 2つの頻度計数、 A 上の F と B 上の G から、 $A \times B$ 上の (直積) 頻度計数 $F \times G$ が次のように生成できる。 $F \times G$ は、 F 内の粒子 (a, n) と G 内の粒子 (b, m) の全ての組合せについて $((a, b), nm)$ という粒子を持つ $(A \times B) \times N$ の部分集合である。これは直積確率測度に対応する。

10

【0063】

[FC II] 写像 $f: A \rightarrow B$ があるとき、頻度計数の写像 $f_*: Freq(A) \rightarrow Freq(B)$ が次のように定義される。頻度計数 F に対して、 $f_*(F)$ は、 $b = f(a)$ なる粒子 (a, m) が少なくとも1つ F 内に存在し n はそのような粒子 (a, m) 全てについての m の和であるような粒子 (b, n) からなる。言い換えれば、集合 $f_*(F)$ は、 F 内の全ての (a, m) について $(f(a), m)$ を追加し、その後、同じ第一成分を持つ異なる粒子がなくなるまで、同じ b の (b, i) と (b, j) を $(b, i+j)$ で置き換えてゆくことで作られる。これは誘導された確率測度に対応する。

20

【0064】

[FC III] もし $A \rightarrow B$ ならば $Freq(A) \rightarrow Freq(B)$ である。つまり、 B 上の頻度計数は自動的に A 上の頻度計数である。 $A \rightarrow B$ で F が A 上の頻度計数であるとき、 F の B への制限 $F|_B$ とは、 F の粒子 (a, n) で a が B に含まれるもの全てからなる B 上の (従って A 上の) 頻度計数である。

【0065】

[FC IV] A 上の2つの頻度計数 F と G は、ある数 $m > 0$ があって、 A の全ての a について $count_F(a) = m \cdot count_G(a)$ であるとき、同値であるといわれる。もし F と G が同値なら、様々な性質がある。 $mass(F) = m \cdot mass(G)$ 、 $supp(F) = supp(G)$ 、 A の全ての a について $P_F(a) = P_G(a)$ 、そして $H(F) = H(G)$ 。

30

【0066】

[FC V] 集合 A について、 A 上の標準頻度計数 $St(A)$ は、 A 内の各 a について粒子 $(a, 1)$ を持つ $A \times N$ の部分集合として定義される。この定義と [FC I] によれば、 $St(A) \times St(B)$ は $St(A \times B)$ と同一であることに注意せよ。

原始写像

【0067】

[PM I] 以下に列挙された全ての原始写像が、原始写像の集合に含まれる。

誘導データと写像

【0068】

ロードされたデータと原始写像に基づいて、そのデータの特徴付ける様々な集合の可能性を探索するために、他のデータと写像が生成される。始めは、入力データが集合上の頻度計数として表現されたものがある。従ってシステムはその集合に適用可能な写像を試しに適用することから始める。そのような写像を適用した結果は新しいデータである。具体的には、プロセスは次のようなデータ構造を維持する：

40

- ・頻度計数の表現を格納するデータ構造 FC 。それは始め、頻度計数として表現された入力データと、入力データがその上にあるような集合の構成要素として現れる全ての集合 A について、その標準頻度計数 $St(A)$ ([FC V] 参照) を持つ。(つまり、もし入力データが $A \times (B \rightarrow C)$ 上の頻度計数なら、 $A, B, C, B \rightarrow C, A \times (B \rightarrow C)$ 上の標準頻度計数が FC 内に含まれるであろう。)それはまた $bool$ や $unit$ などのいくつかの標準的集合上の標準頻度計数も含む。

50

・集合の記号表現を格納するデータ構造 $S E T S$ 。始めそれは $F C$ 内の頻度計数がその上にあるような集合を含む。

・写像の記号表現を格納するデータ構造 $M A P S$ 。始めそれは原始写像を含む。

【0069】

過程が続く上で、 $F C$ 、 $S E T S$ 、 $M A P S$ に、以下の何れかの様に要素が加えられる。

【0070】

[$D I$] もし $F C$ 内に頻度計数の組 F 、 G があれば、 $F \times G$ が $F C$ に加えられ得る ([$F C I$] 参照)。3つ以上の頻度計数の組についても同様である。

【0071】

[$D I I$] もし $M A P S$ 内の写像が $M A P S$ 内の写像に適用できるなら (例えば [$P M I I I$]、[$P M I V$]、[$P M V$]、[$P M V I$]、[$P M X I I$])、適用した結果の写像が $M A P S$ に加えられ得る。例えば、いくつかの写像の組が選ばれてそれらの積写像あるいは可能ならそれらの合成が $M A P S$ に加えられ得る。あるいは任意の写像が他の写像に適用されて結果が $M A P S$ に加えられ得る。

10

【0072】

[$D I I I$] $S E T S$ 内の集合の部分集合が $S E T S$ に加えられ得る。頻度計数が部分集合に制限されうる。部分集合の逆像が $S E T S$ に加えられ得る。 A の部分集合 B について、部分集合判別写像 $subset_B : A \rightarrow \{true, false\}$ ($a \in B$ なら $subset_B(a) = true$ 、それ以外なら $false$ と定義される) が $M A P S$ に加えられ得る。

【0073】

20

[$D I V$] もし集合 A 上の頻度計数 F が $F C$ 内に、写像 $f : A \rightarrow B$ が $M A P S$ 内にあれば、 $f_*(F)$ が $F C$ に加えられ得る ([$F C I I$] 参照)。このルールを使って頻度計数が $F C$ に加えられるときは、 $F C$ は使われた写像も記録する。

【0074】

集合を頂点とし、写像を辺として、これらの集合は有向グラフ構造を形成すると考えられる。集合上の頻度計数もまた、頻度計数を頂点とし、写像を辺として、有向グラフ構造を形成すると考えられる。

【0075】

これらの写像とデータは様々な順番でこれらのデータ構造に加えることができる。例えば、上記の木構造内で幅優先探索の順序を使うことができる。本実施例では、確率的アルゴリズムが使われる：

30

【0076】

探索アルゴリズム

【0077】

概要

【0078】

以下の1から6の動作の何れかを確率的に実行せよ：

1. $F C$ 内の頻度計数 F と G の組を選んで $F \times G$ を $F C$ に加える。 F が集合 A 上、 G が集合 B 上であるとして、 $A \times B$ を $S E T S$ に加える。

2. [$D I I$] に従って写像に適用可能な $M A P S$ 内の写像を選んで適用し、結果を $M A P S$ に加える。

40

3. $S E T S$ 内の集合 A を選び、 A のある真部分集合 B を $S E T S$ に加え、 $subset_B : A \rightarrow \{true, false\}$ を $M A P S$ に加える。

4. $F C$ 内の頻度計数 F を選ぶ。 F が集合 A 上であるとして、 $S E T S$ 内の A の真部分集合 B を選び、 $F|_B$ を $F C$ に加える。

5. $M A P S$ 内の写像 $f : A \rightarrow B$ を選び、 $S E T S$ 内の B の真部分集合 C を選ぶ。逆像 $f^{-1}(C)$ を $S E T S$ に加える。

6. $F C$ 内の頻度計数 F と、 F が上にある集合から何か他の集合への写像 f を $M A P S$ 内に選び、 $f_*(F)$ を $F C$ に加える。

【0079】

50

詳細

【0080】

図2は探索アルゴリズムのフローチャートを示す。動作とその対象の選択は確率的になされる。

【0081】

【0104】

原始写像

【0105】

一般的な原始写像に加えて、画像に特に有用な原始写像を加えることもある。例えば、画像が普通そうであるように画素からなれば、画素間の隣接関係が有用であるかもしれない。これはDomの2要素が隣接画素であるときのみtrueを返す原始写像

$Nb : Dom \times Dom \rightarrow bool$ としてシステムに入れることができる。もう一つの例は、例えばウェーブレットフィルターなどの、画像処理の関連分野で知られる種々のフィルターである。

【0106】

誘導データと写像

【0107】

本方法がMAPSとFCに加えそうな写像とデータの簡単な例を挙げる：

【0108】

A．色頻度

1．A1．[D I]より、2つの頻度計数 $Dom \times Col$ 上のImとDom上のSt(Dom)に基づいて、 $(Dom \times Col) \times Dom$ 上の頻度計数 $Im \times St(Dom)$ がFCに加えられる。

2．A2．[D I V]より、A1からの $Im \times St(Dom)$ と(原始写像なのでMAPS内にある)評価写像 $ev : (Dom \times Col) \times Dom \times Col$ に基づいて、 $ev * (Im \times St(Dom))$ がFCに加えられる。 Col 上の頻度計数 $ev * (Im \times St(Dom))$ は粒子 (c, n_c) の集合で、 n_c は色cを持つ画素の数である。

【0109】

B．色の違いと位置の違いの頻度

1．B1．[D I I]により、対角写像 $diag : (Dom \times Col) \times (Dom \times Col) \times (Dom \times Col)$ 、積写像 $mp : (Dom \times Col) \times (Dom \times Col) \times (Dom \times Dom \times Col \times Col)$ 、及び対角写像

$diag : Dom \times Dom \times (Dom \times Dom) \times (Dom \times Dom)$ に基づいて、写像 $(mp \circ diag) \times diag : (Dom \times Col) \times (Dom \times Dom) \times (Dom \times Dom \times Col \times Col) \times (Dom \times Dom) \times (Dom \times Dom)$ がMAPSに加えられる。

2．B2．[D I I]より、評価写像 $ev : (Dom \times Dom \times Col \times Col) \times (Dom \times Dom) \times Col \times Col$ 、及び $Dom \times Dom$ の恒等写像に基づいて、写像 $ev \times id_{Dom \times Dom} : (Dom \times Dom \times Col \times Col) \times (Dom \times Dom) \times (Dom \times Dom) \times (Col \times Col) \times (Dom \times Dom)$ がMAPSに加えられる。

3．B3．[D I I]より、色空間の引算と画像領域の差写像に基づいて、写像 $Sub_{Col} \times Diff_{Dom} : (Col \times Col) \times (Dom \times Dom) \times Col \times V_{Dom}$ がMAPSに加えられる。

4．B4．[D I I]によりB1、B2、B3でMAPSに加えられた写像を合成して $(Sub_{Col} \times Diff_{Dom}) \circ (ev \times id_{Dom \times Dom}) \circ ((mp \circ diag) \times diag) : (Dom \times Col) \times (Dom \times Dom) \times Col \times V_{Dom}$ がMAPSに加えられる。

5．B5．[D I]により、 $(Dom \times Col) \times (Dom \times Dom)$ 上の頻度計数 $Im \times St(Dom \times Dom)$ がFCに加えられる。

10

20

30

40

50

6. B6. [D I V]により、B4の写像をB5で加えられた頻度計数 $I m \times S t (D o m \times D o m)$ に適用した結果がFCに加えられる。

B6で加えられた $C o l \times V_{D o m}$ 上の頻度計数は粒子 $((d,), n_d,)$ の集合で、 $n_d,)$ は i) 色の違い d を持ち、 $i i$) それらの間の画像領域内でのベクトルである、画素の組の数である。

【0110】

パターン

【0111】

A2で得られる $C o l$ 上の頻度計数 $e v_* (I m \times S t (D o m))$ は、あまり多くの色が使われていないときに小さいエントロピーを持つ。画像全体が一色なら、エントロピーとして可能な最小値0を持つ。

10

【0112】

B6で加えられる $C o l \times V_{D o m}$ 上の頻度計数は、同じ特定の色の違いと同じベクトルで隔てられたがその組がたくさんあるとき、小さいエントロピーを持つ。例えばもし一つの色の水平な直線があれば、色の違い0と水平ベクトルの粒子に比較的高い集中(計数の高い粒子)があり、この頻度計数のエントロピーは低くなる。

例2: データマトリクス

【0113】

データマトリクスとはN行D列の直方配列で、各行が異なる観察あるいは固体を与え、各列が異なる属性や変数を与えるものである。各変数は、ここで値集合と呼ぶ何かの集合の要素である値をとることができる。例えば、もし変数が制すうちだけを取り得るなら、値集合は整数の集合である。もし変数が任意の数を取り得るなら、値集合は実数の集合である。あるいは変数が「yes」か「no」の値しかとれないなら、値集合はブール値の集合でありうる。

20

【0114】

D個の変数を a_1, a_2, \dots, a_D で表し、それらの変数が値を取り得る集合をそれぞれ X_1, X_2, \dots, X_D で表す。すると、各観察は集合 $X_1 \times X_2 \times \dots \times X_D$ の要素を与える。データマトリクスの形をとった入力データは、本実施例では各観察が一つの粒子の1計数に寄与する $X_1 \times X_2 \times \dots \times X_D$ 上の頻度計数として表される。従って頻度計数のマスはNである。

30

【産業上の利用可能性】

【0115】

以上のように、高次元データがより効果的に解析され得るように、また与えられたデータ内によりよくパターンを見つけられるように、与えられたデータを配置する方法と装置が開示された。本発明は広い範囲の産業で利用可能である。それらの産業では、ますます多くのデータが収集され、巨大なデータの山から適切な情報を見つけ出すことがますます重要になっている。本発明が有用な分野は、多数の遺伝子と、ある遺伝病に罹患する比較的少数の患者の例、及び用意に百万次元(画素)を持ちうる画像の例を含む。

【0116】

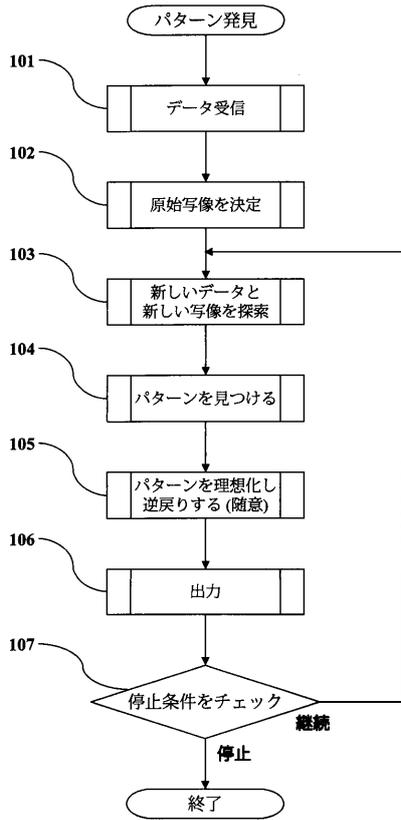
本明細書には本発明の特定の好ましい特徴のみを説明記述したが、関連技術分野で熟練した者には多くの修正や改変が思い浮かぶであろう。例えば、本発明を説明するためにここで使われた集合や写像の概念は、様々な分野で多くの同値あるいは類似の概念を持つ。例えば、関数、型、メソッド等である。集合や写像などの用語は、望むならば完全に避けることが可能である。本発明全体をデータとサブルーチンの言葉で記述することも可能である。しかし、そのような表面的な違いは、真の違いではない。

40

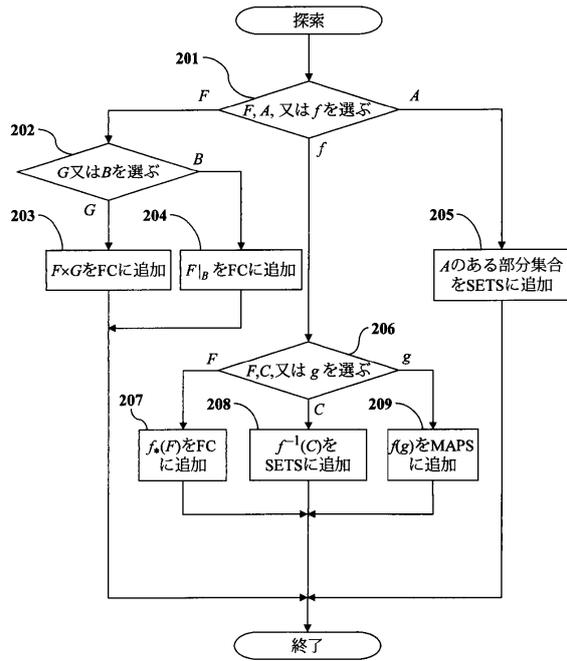
【0117】

従って、添付の特許請求の範囲は、そのような全ての修正、変更、用語の違いを本発明の真の精神のうちに入るものとして全て含むよう意図されたものであることを理解されたい。

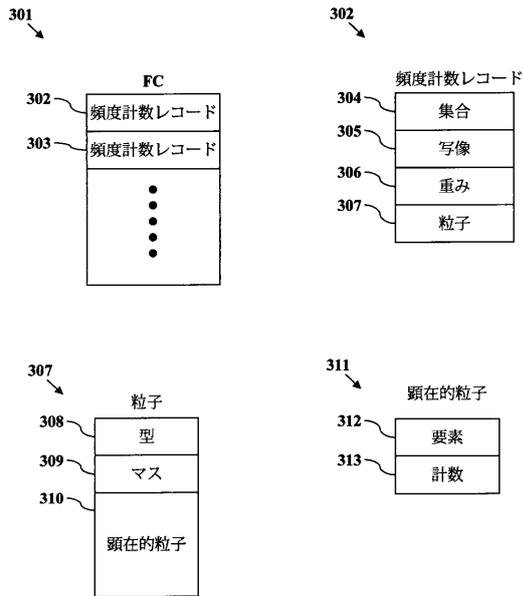
【図1】



【図2】



【図3】



【図4】

