

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4916614号
(P4916614)

(45) 発行日 平成24年4月18日 (2012. 4. 18)

(24) 登録日 平成24年2月3日 (2012. 2. 3)

(51) Int. Cl.

F I

G 0 6 N 3 / 0 0 (2006.01)

G 0 6 N 3 / 0 0 5 5 0 C

請求項の数 1 (全 60 頁)

(21) 出願番号	特願2000-615965 (P2000-615965)	(73) 特許権者	390023674
(86) (22) 出願日	平成12年4月19日 (2000. 4. 19)		イー・アイ・デュポン・ドウ・ヌムール・
(65) 公表番号	特表2002-543538 (P2002-543538A)		アンド・カンパニー
(43) 公表日	平成14年12月17日 (2002. 12. 17)		E. I. DU PONT DE NEMO
(86) 国際出願番号	PCT/US2000/010425		URS AND COMPANY
(87) 国際公開番号	W02000/067200		アメリカ合衆国、デラウェア州、ウイルミ
(87) 国際公開日	平成12年11月9日 (2000. 11. 9)		ントン、マーケット・ストリート 100
審査請求日	平成19年4月18日 (2007. 4. 18)		7
(31) 優先権主張番号	60/131, 804	(74) 代理人	100127926
(32) 優先日	平成11年4月30日 (1999. 4. 30)		弁理士 結田 純次
(33) 優先権主張国	米国 (US)	(74) 代理人	100140132
(31) 優先権主張番号	09/466, 041		弁理士 竹林 則幸
(32) 優先日	平成11年12月17日 (1999. 12. 17)		
(33) 優先権主張国	米国 (US)		
前置審査		最終頁に続く	

(54) 【発明の名称】 実験データの分布状階層的発展型モデリングと可視化の方法

(57) 【特許請求の範囲】

【請求項 1】

増幅 DNA フラグメントの同定に関連するグローバルな情報コンテンツを有するフィーチャー集合を選択する、コンピューターにより実行される方法であって、多数の入力データ点と対応する出力データ点が取得されてデータ集合が規定され、該取得された入力及び出力データ点が記憶装置内に記憶され、該フィーチャー集合およびデータ集合は DNA 増幅プロセスの間に生成される方法に於いて、

(a) 複数のフィーチャー部分空間を創るが、各前記フィーチャー部分空間が該データ集合からのフィーチャー集合を含むように、該創る過程と、

(b) 該データ集合の該入力を量子化するが、該入力が値の範囲を有し、それは該値の範囲を部分範囲に分け、それにより前記フィーチャー部分空間を複数のセルに分けることによりするよう、該量子化する過程と、

(c) 少なくとも 1 つのローカルセルのニシのエントロピー E を計算し、ニシのエントロピー E の補数としてローカルエントロピー W ($W = 1 - E$) を規定し、そして各セルのローカルエントロピー W の加重和を計算することによって、各フィーチャー部分空間のグローバルエントロピーを計算する過程と、

(d) 該グローバルエントロピーに基づき、グローバルな情報コンテンツを有する少なくとも 1 つのフィーチャー集合を選択する過程とを具備する、方法。

【発明の詳細な説明】

【0001】

10

20

【発明の分野】

本発明は " 対象 (objects) " の階層 (hierarchy)、例えば、フィーチャー (features)、モデル (models)、フレームワーク (frameworks)、そしてスーパーフレームワーク (super-frameworks)、を創るために、データの画像的表現の概念を情報理論 (information theory) からの概念と組み合わせる。本発明はシステムの実験型モデルを、前に取得されたデータ、すなわち、該システムへの入力と該システムからの対応する出力を表すデータ、に基づいて創る方法と機械可読記憶媒体 (machine readable storage medium) とに関する。次いで該モデルは次の取得入力からシステム出力を正確に予測するため使われる。本発明の方法と機械可読記憶媒体は情報理論と熱力学の原理に基づく、エントロピー関数を使用し、該方法は複雑な、多元処理 (multi-dimensional process) のモデリングに特に好適である。本発明の方法はカテゴリー的モデリング (categorical modeling)、すなわち、出力変数が離散的状態 (discrete states) をとる場合、及び定量的モデリング、すなわち、出力変数が連続的な場合、の両者に使用出来る。本発明の方法は、外見には混乱したシステムであるように見えるものの下にある順序、又は構造を顕わすために、データ集合の最適表現、すなわち最も情報豊富な表現 (most information-rich representation) を同定 (identifies) する。発展型プログラミング (evolutionary programming) の使用は最適表現を同定する 1 方法である。該方法は多元的フィーチャー空間 (multi-dimensional feature spaces) の情報コンテンツ (information content) を特徴付ける中でローカル及びグローバルの両情報メジャー (both local and global information measure) のその使用により際だっている。実験はローカル情報メジャーがモデルの予測能力 (predictive capability) を支配することを示した。かくして、全体のデータ集合上でのグローバルな最適化を主として使う、多くの他の方法と対照的に、本方法はグローバルに影響されるが、ローカルに最適化される技術、として説明出来る。

【 0 0 0 2 】

【発明の技術的背景】

情報理論

システムの情報コンテンツを説明するためにエントロピー関数 (entropy function) を使用する思想は、彼のパイオニア的業績、1948年発行の、ベルシステムテクニカルジャーナル (Bell System Technical Journal)、27, 379 - 423, 623 - 656、" 通信の数学的理論 (A Mathematical Theory of Communication) " でシー・イー・シャノン (C. E. Shannon) により初めて導入された。シャノンは統計力学での対応する定義と形式的に同様なエントロピーの定義が起こり得るイベントの総体 (ensemble) 内での特定のイベントの選択から得られる情報を測定するため使用出来ることを示した。シャノンのエントロピー関数は下記で表され、

【 0 0 0 3 】

【数 1】

$$H(p_1, p_2, \dots, p_n) = - \sum_{k=1}^n p_k \ln p_k$$

【 0 0 0 4 】

ここで p_k は第 k 番目のイベントの発生確率を示し、ユニークに下記 3 条件を満足する、
 1. $H(p_1, \dots, p_n)$ は $k = 1, \dots, n$ で $p_k = 1/n$ で最大となる。これは均一な確率分布が最大エントロピーを有することを意味する。加えて、 $H_{\max}(1/n, 1/n, \dots, 1/n) = \ln n$ 。従って、均一確率分布のエントロピーは起こり得る状態の数と共に対数的に縮尺 (scales) する。

2. $H(AB) = H(A) + H_A(B)$ ここで A と B は 2 つの有限スキーム (finite schemes) である。 $H(AB)$ はスキーム A と B の全エントロピーを表し、 $H_A(B)$ はスキーム B を与えられたスキーム A の条件的エントロピーである。該 2 つのスキーム分布が相互に独立の時、 $H_A(B) = H(B)$ である。

3. $H(p_1, p_2, \dots, p_n, 0) = H(p_1, p_2, \dots, p_n)$ 。スキーム内の発生確率ゼロのどんなイベントもエントロピー関数を変化させない。

【0005】

シャノンの仕事は1次元の電気信号の情報コンテンツを説明することに向けられた。1998年に、ケンブリッジ大学プレス (Cambridge University Press) で発行された彼の本、フィッシャー情報からの物理学：ユニファイケーション (Physics from Fisher Information: A Unification) で、ロイフリーデン (Roy Frieden) は "シャノンエントロピー (Shannon Entropy)" を全体のデータ集合間のグローバルな情報メジャーとして説明している。"フィッシャーエントロピー (Fisher entropy)" として知られる、代替りの情報メジャーも又データ集合間のローカルな情報の測定量としてフリーデンにより説明されている。数学的モデル化で、フリーデンはフィッシャーエントロピーが物理的法則を発見するために特に好適であることを最近示した。

10

【0006】

より最近に、テー・ニシ (T. Nishi) はどんなデータ集合にも適用出来る、正規化された "情報エントロピー" 関数を規定するために該シャノンのエントロピー関数を使用した。1991年、京都、325、材料の機械的挙動に関する国際会議論文集 (Proceedings of the International Conference on 'Mechanical Behaviour of Materials VI')、ハヤシ、テー・及びニシ、テー・ (Hayashi, T. and Nishi, T.) 著、"ポリマーアロイの形態学と物理的特性 (Morphology and Physical Properties of Polymer Alloys)"、参照。1992年発行、高分子論文集 (Kobunshi Ronbunshu)、49(4)、373-82、ハヤシ、テー・、ワタナベ、エイ・、タナカ、エイチ・及びニシ、テー・ (Hayashi, T., Watanabe, A., Tanaka, H. and Nishi, T.) 著、"3成分不相溶性ポリマーアロイの形態学と物理的特性 (Morphology and Physical Properties of Three-Components Incompatible Polymer Alloys)" 参照。

20

【0007】

ニシの定義は次ぎの様に抄録されるが、 n のデータ要素 (data elements) を有するデータ集合 (data set) $D = \{d_1, \dots, d_n\}$ を考える。もし全要素の和 d_{tot} が次の様に定義されるならば、

【0008】

【数2】

30

$$d_{tot} = \sum_{i=1}^n d_i,$$

【0009】

d_{tot} は、

【0010】

【数3】

$$f_i = d_i / d_{tot} \quad \forall i \in \{1, \dots, n\}.$$

40

【0011】

の様に該データ要素の各々を正規化 (normalize) するため使用出来る。

次いで、情報エントロピー関数 (informational entropy function)、 E を次の様に規定することが出来る、

【0012】

【数4】

$$E = (\sum_i f_i \ln f_i) / \ln(1/n).$$

50

【 0 0 1 3 】

該エントロピー関数 E はそれが 0 と 1 の間に正規化される有用な特性 (property) を有する。 $f_i = 1 / n$ の、完全に均一な分布 (perfectly uniform distribution) は 1 の E 値となる。該分布がより不均一になるにつれ、 E の値は低下し漸近的にゼロに近づく。該ニシの情報エントロピー関数 E の顕著な利点はそれが分布の形状に無関係にどんな分布の均一性も特徴付けることである。対照的に、普通使用される "標準偏差 (standard deviation)" はガウス分布 (Gaussian distribution) 用でのみ標準的統計 (standard distribution) に入ると通常解釈される。

【 0 0 1 4 】

ニューラルネットワーク (neural networks)、統計的回帰 (statistical regression)、決定木法 (decision tree methods) の様な従来技術の方法は或る本質的限定を有する。ニューラルネットワークと他の統計的回帰方法はカテゴリー的モデリングに使用されて来たが、それらは、該ネットワークのノード内で使用される連続非線形シグモイド関数 (continuous non-linear sigmoid function) のために、定量的モデル化に遙かにより適合し、より良く動作する。決定木は、連続的出力値に関する精確な定量的予測をする能力に欠けるためにカテゴリー的モデリングに最も良く適合している。

【 0 0 1 5 】

【発明の概要】

本発明は情報エントロピーの概念を一般化し、それらの概念を多次元データ集合へ延長している。特に、シャノンにより表明された情報エントロピーの定量化は修正され、1つ以上の入力、又はフィーチャー、と1つ以上の出力とを有するシステムから得られたデータに適用される。情報豊富 (information-rich) でありかくして該システム出力 (含む複数) の予測に有用なデータ入力の種々の部分集合 (subset)、又はフィーチャーの部分集合を同定 (identify) するためにエントロピー定量化 (entropy quantification) が行われる。又該エントロピー定量化は情報豊富な種々のフィーチャー部分集合内で領域 (region)、又はセル (cell) を同定する。該セルは固定的又は適合的なビンニング過程 (binning process) を使用してフィーチャー部分空間内で規定される。

【 0 0 1 6 】

入力組み合わせ (input combination)、又は特徴組み合わせ (feature combination)、はフィーチャー部分空間を規定する。該フィーチャー部分空間は2進ビット記号列 (binary bit string) により表され、ここでは遺伝子 (genes) として引用される。遺伝子はどの入力 that 特定部分空間にあるかを示し、従って特定の部分空間の次元数 (dimensionality) は該遺伝子数列 (genes sequence) の "1" のビットの数により決定される。望ましい情報特性を有する部分空間に対応するそれら遺伝子を同定するために全てのフィーチャー部分空間の情報豊富さがエグゾースチブ (exhaustively) に探索される。

【 0 0 1 7 】

起こり得る部分空間 (possible subspace) の全数が少なければ、エグゾースチブな探索が最も情報豊富な部分空間を同定する好ましい方法であることは注意すべきである。多くの場合、しかしながら、起こり得る部分空間の数は全ての起こり得る部分空間をエグゾースチブに探索することが計算的に非現実的である程充分大きい。それらの状況では、該部分空間は遺伝子数列を操作する遺伝的アルゴリズムを使用して探索されるのが好ましい。すなわち、遺伝子は望ましい情報特性を有するフィーチャー部分空間の集合を進化させるよう組み合わせられ及び/又は選択的に突然変異 (mutated) させられる。特に、該遺伝的フィーチャー部分空間進化過程 (evolution process) 用の適応度関数 (fitness function) はその特定の遺伝子により表されるフィーチャー部分空間用情報エントロピーのメジャー (measure) である。情報コンテンツの他のメジャーは該出力に関する該部分空間の均一度を示す (measure)。これらのメジャーは分散 (variance)、標準偏差、又は或るしきい値を越える指定出力依存確率を有するセルの数 (又はセルのパーセンテージ) の様な発見的方法 (heuristics) を含む。これらの情報のメジャーは望ましい情報特性、すなわち高い情報コンテンツを有する遺伝子、又は部分空間を同定するために使用されてもよい。加え

10

20

30

40

50

て、決定木ベースの方法が使用されてもよい。これらの代替えの方法はエグゾースチブな探索を行う時望ましい部分空間を同定するため使用されてもよい。

【 0 0 1 8 】

好ましい実施例では、ここではグローバルエントロピーと呼ぶ、該フィーチャー部分空間エントロピーは、該部分空間内のセルのエントロピーメジャーの加重平均を計算することにより決定されるのが好ましい。出力特定のエントロピーメジャーも又使用されてもよい。セルエントロピーはここではローカルエントロピーと呼ばれ、修正されたニシのエントロピー計算を使用して計算される。

【 0 0 1 9 】

実験型モデルが次いで階層的な仕方で創られるが、それは、高い情報コンテンツを有するよう決定されたフィーチャー部分空間の組み合わせを調べることによる。フィーチャー部分空間は、テストデータ（既知の対応出力を有するサンプル入力データ点）を使用する高精度の予測を提供するフィーチャー部分空間の組み合わせを見出すためにエグゾースチブな探索技術を使用して選択されそしてモデル内へ組み合わせられる。該モデルは又遺伝的アルゴリズムを使用して発展させられてもよい。この場合、該モデル遺伝子はどのフィーチャー部分空間が使用されるかを指定し、該モデル遺伝子の長さは望ましい情報特性を有するとして前に同定されたフィーチャー部分空間の数により決定される。該モデル発展過程で使用される該適応度関数は考慮下の特定モデルの予測精度であるのが好ましい。

【 0 0 2 0 】

本発明の 1 側面に依れば、次ぎに取得される入力からシステム出力を精密に予測するため、該システムへの対応する入出力を表す、前に取得されたデータに基づきシステムの実験型モデルを創る方法が提供される。該方法は、

（ a ）該システムへの多数の入力と対応する該システムからの出力とからデータ集合を取得する過程と、

（ b ）該前に取得したデータ集合を、少なくとも 1 つのトレーニングデータ（training data）集合と、少なくとも 1 つのテストデータ（test data）集合と、そして少なくとも 1 つの検証データ（verification data）集合とにグループ分けする過程を具備しており、該集合は相互に一致してもよく、或いは前に取得したデータの排他的（exclusive）又は非排他的（non-exclusive）部分集合であってもよく、該方法は又、

（ c ）高いグローバルエントロピー加重（weights）を有する複数のフィーチャー部分空間を、

（ i ）前記トレーニングデータ集合からフィーチャー部分空間を規定する複数の入力を選択する過程と、

（ i i ）固定的か又は適合的か何れかの量子化方法（quantization）により、各入力範囲を部分範囲（subrange）に分けることにより該フィーチャー部分空間をセルに分ける過程と、

（ i i i ）ローカルセラーエントロピー加重による加重平均か、又は出力特定のエントロピー加重による加重平均か何れかを形成することにより、グローバルエントロピー加重を決定する過程と、

により決定する過程と、

（ d ）オプション的に、高いエントロピー加重を有する該決定されたフィーチャー部分空間内での各入力発生頻度を調べ、削減された次元数データ集合を規定するために最も頻繁に発生するそれらの入力のみを保持する過程と、そしてその後過程（ c ）を繰り返す過程と、

（ e ）オプション的に、該削減された次元数フィーチャーデータ集合を規定するようにシステム入力から最も精密にシステム出力を予測する最適又は最適に近い次元数と最適又は最適に近い量子化条件を決定するために、複数の量子化条件下で該削減された次元数データ集合の複数の該次元（例えば、該次元の幾つか、又は全て）上でエグゾースチブに探索する過程と、

（ f ）前記データ集合上のシステム入力からシステム出力を最も精密に予測する高いグロ

10

20

30

40

50

ーバルエントロピー加重（例えば、フィーチャーデータ集合の部分か、又は全体か何れか）を有する該決定されたフィーチャー部分集合の組み合わせを決定する過程と、

（g）テストデータ集合上でシステム入力からシステム出力を最も精密に予測する削減された次元数のフィーチャーデータ集合に部分集合（例えば、削減された次元数のフィーチャーデータ集合の部分か、又は全体かの何れか）を決定する過程とを具備している。

【0021】

大きなデータ集合用には、該モデル創生過程（b）-（g）は、次いで最適モデルのグループを見出すために種々のトレーニング及びテストデータ集合上で繰り返されてもよい。この最適モデルのグループはそれらのモデルから生じる1つ以上の予測を開発するために新しいデータについて“ポール（polled）”されてもよい。これらの予測は、例えば、勝者1人占め（winner-takes-all）の投票ルールに基づいてもよい。システム入力から最も精密にシステム出力を予測する最適モデルのグループの部分集合は次いで次の様に決定される。テストデータ集合の入力がモデルの選択された部分集合のグループの各モデルに從属させられ（ランダムに選択されてよい）、各部分集合で予測された出力は各テストデータ出力と比較される。該部分集合で予測された出力の計算過程は（b）-（e）{又はオプションとして（b）-（g）}と同様な仕方で行われ、そこでは個別のモデル出力予測値を入力として、実際の出力値を出力として使用して新しいトレーニング及びテストデータ集合が創られる。この過程はモデルの多数の選択された部分集合グループ用に繰り返されてもよい。モデルの該選択された部分集合グループは次いで、“フレームワーク”を規定するためにシステム入力からシステム出力を最も精密に予測するモデルの最適部分集合グループを見出すために発展（evolved）させられる。

【0022】

フレームワーク創生過程は、最適フレームワークのグループを見出すために、モデル創生過程と同様な仕方でも更に繰り返されてもよい。最適フレームワークのこのグループは、それらのフレームワークから生じる1つ以上の予測を開発するために新データ上で“ポール”され得る。これらの予測は、例えば、勝者1人占めの投票ルールに基づくことが出来る。システム入力からシステム出力を最も精密に予測する最適フレームワークのグループの部分集合は次いで次の様に決定される。テストデータ集合の入力はフレームワークの該選択された部分集合グループの各フレームワークに印加され、各フレームワーク部分集合で予測された出力が各テストデータ出力と比較される。該部分集合で予測される出力の計算過程は（b）-（g）と同様な仕方で行われ、そこでは個別モデルフレームワークで予測された値を入力としてそして実際の出力を出力として使用して新トレーニング及びテストデータ集合が創られる。この過程はフレームワークの多数の選択された部分集合グループ用に繰り返される。フレームワークの該選択された部分集合グループはシステム入力からシステム出力を最も精密に予測する、“スーパーフレームワーク”と呼ばれる、フレームワークの最適部分集合グループを見出すために発展させられる。

【0023】

最適モデル決定過程、最適フレームワーク決定過程、又は最適スーパーフレームワーク決定過程は予め決められた停止条件が達成されるまで繰り返される。該停止条件は、例えば、1）発展型対象の族（family of evolutionary objects）のポーリングから予め決められた予測精度の達成、又は2）予測精度でのインクレメンタルな改善が予め決められたしきい値より低下した時、又は3）予測精度での更に進んだ改善が達成されない時、として規定されてもよい。

【0024】

分布状階層的発展（Distributed hierarchical evolution）は、モデル、フレームワーク、スーパーフレームワーク他の様な逐次的により複雑に相互作用する発展型“対象”のグループが、逐次的により大量の複雑なデータをモデル化し理解するために、創られる発展型の過程である。

【0025】

【本発明の詳細な説明】

10

20

30

40

50

図 1 は本発明の方法 1 0 0 の全体的流れを図解するブロック線図である。この図から評価される様に、実験データから複雑なシステムのモデルを創生するために発展型過程 (evolutionary process) が使用される。好ましい方法は、" 発展型対象 (evolutionary objects) "、例えば、フィーチャー 1 3 0、モデル 1 4 0、フレームワーク 1 5 0、そしてスーパーフレームワーク 1 6 0 他、の伸展する階層 (extensible hierarchy) を創るために、データ 1 1 0 の多次元的表現を情報理論 1 2 0 と組み合わせる。該過程は 1 7 0 で示した階層的な仕方です更に組み合わせを発生するため続けられ得る。

【 0 0 2 6 】

最初に、フィーチャー部分空間 (feature subspace) と呼ばれる、入力 of の組み合わせは、初期のランダムに選択されたフィーチャー部分空間プールからエグゾースチブな探索 (exhaustive search) 又は発展型の過程により、同定 (identified) される。次いでモデルを創るためにフィーチャー部分空間の最適組み合わせ (optimum combination) が探索されるか又は発展 (evolved) させられ、フレームワークを創るためにモデルの最適組み合わせが更に探索されるか又は発展させられ、そしてスーパーフレームワーク他を創るためにフレームワークの最適組み合わせが更に探索されるか又は発展させられる。上記説明のより複雑な発展型対象の逐次的発展は、予め決められた停止条件、例えば、予め決められたモデル性能、が達成されるまで続く。ルールとして、該データ集合 (data set) が大きい程、これらの対象のより多くが創られるので、実験型モデル (empirical model) の複雑さは、該入力の、該データが取得された該システムの出力との相互作用の複雑さを反映する。

【 0 0 2 7 】

ここに説明した方法の展開で、幾つかの設計基準 (design criteria) が考えられた。該方法が、任意の非線形構造を有するデータ空間 (data space) を成功裡に処理することが必要である。該方法が、入力を知って出力を予測する " 前向き (forward) " 問題と、出力を知って入力を予測する " 逆向き (inverse) " 問題との間を区別せず、それによりデータのモデル化と制御の問題を同じ足場 (footing) 上に置くことも又望ましい。これは該データ集合それ自身の上に最小の追加的モデルジオメトリー (additional model geometry) だけが重ね合わされることを意味する。用語 " ジオメトリー (geometry) " は、回帰技術 (regression technique) で導入される様な、線形及び非線形の両多様性を含む。対称性 (symmetry) もここでは目下のモデリングタスク用に最も情報豊富な (information-rich) 入力又は入力の組み合わせを同定する利点を有する。この知識は意志決定及び計画用の最適戦略を開発するため使用され得る。最後に、該方法は、それが事実便利に実施されるために計算的に扱い易い (tractable) 必要がある。これらの設計目標を充たすために、幾つかの現在の線形及び非線形な方法が注意深く解析され、共通のテーマが基本的な限定と機会とを同定する目標を用いて要約された。

【 0 0 2 8 】

下記の議論は情報理論及び発展からの概念を使用して 1 つのモデルの発展の基本的方法を説明することから始まる。より大きい。より複雑なデータ集合を説明するために逐次的により複雑な対象の逐次的で階層的な発展に向かうために該方法を更に伸展させることが次ぎに説明される。データ出力がなくても入力フィーチャークラスター (input feature cluster) を発見する方法の下にある原理の応用が次いで論じられ、それに多次元データ空間内で " 情報可視化 (information visualization) " を行う方法の説明が続く。ハイブリッドのモデリングスキームを創るために本発明の方法をニューラルネットワーク (neural networks) の様な他のモデリングパラダイム (modeling paradigms) と組み合わせることが次いで詳述される。該説明は、遺伝的プログラミング (genetic programming) の分野と結合された本発明の方法のデータモデル化の取り組みを使用して物理的法則を発見する、新しい取り組みを結論としている。

【 0 0 2 9 】

関心の点として、情報理論からの基本的アイデアは全てのこれらの問題を解くに必要なコアツール (core tools) を提供し、簡単で統合的核 (simple, unifying kernel) を該方

10

20

30

40

50

法に提供することは述べるに値する。エントロピー (entropy) の概念はデータ空間内の秩序 (order) { 又は混乱 (disorder) } の定量的メジャー (quantitative measure) を提供する。このメジャーは、初期に混乱したシステムからの秩序の発生をドライブする発展型エンジン用の適応度関数 (fitness function) として使用され得る。この意味で、情報理論はドライバーを提供し、発展型プログラミングは発見過程をシステム化するエンジンを提供する。最後に、本発明の方法で説明されるパラダイムはデータドライブされている (is data driven) が、それはデータ自身の中の情報コンテンツ (information content) が予測 (prediction) に使用されるからである。かくして、該方法は、下にある数学のその固有の制限を有する数学的モデル化の分野と反対に、実験型モデル化の分野に真正面 (squarely) から属する。

10

データモデリング (DATA MODELING)

情報エントロピーの概念に基づくフレームワークは、入力 of 集合を与えられたとして1つか又は多数か何れかの出力が予測される必要がある様な、データモデリングの問題に適用されて来た。基本的方法は次の過程から成るが、すなわち

1. データ表現 (data representation) 又はデータ事前処理 (data preprocessing)、
2. セル境界 (cell boundary) を規定する固定的又は適合的 (adaptive) な方法を使用するデータ量子化 (data quantization)、
3. 遺伝的発展及び情報エントロピーを使用するフィーチャー組み合わせ選択、
4. システム入力からシステム出力を最も精密に予測するフィーチャーデータ集合の部分集合 (subset) の決定である。

20

1. データ表現

典型的な実験的に得られたデータ集合で、幾つかの "測定" 入力と出力とが提供される。各システム入力とシステム出力は、ここでデータ点 (data points) と呼ぶ、データ値の入力及び出力のシーケンスを得るようにサンプリングされるか他の仕方で測定される。目標 (goal) は該データ点出力を最も正確に予測するために該データ点入力から最大の情報を抽出することである。多くの実システム (real system) では、該データ点、又は実際の測定された入力は、それらが該データの適切な表現として留まるに十分な程 "情報豊富 (information-rich)" である。他の場合は、これはそうでないかも知れず、該データを表現するより適切な "固有ベクトル (eigenvectors)" を創るために該データを変換することが必要かも知れない。共通に使用される変換には特異値分解法 (singular value decomposition) { エスブイデー (SVD) }、主成分分析法 (principal component analysis) { ピーシーエイ (PCA) }、部分的最小2乗法 (partial least square method) { ピーエルエス (PLS) 法 } が含まれる。

30

【0030】

最も大きい対応する "固有値 (eigenvalues)" を有する主成分 "固有ベクトル" (eigenvectors) が該データモデリング過程用入力として通常使われる。該主成分選択法には2つの顕著な限定がある。

【0031】

a. 該主成分法は入力の分散のみを取り扱い、出力に関する情報は何もエンコードしない。多くのモデリング問題で、モデル化されつつある出力特性に関する最も多くの情報を含む比較的低い固有値を有するのは固有ベクトルである。

40

【0032】

b. 該ピーシーエイ法は入力の線形変換を行う。これは全ての問題用には、特に入力 - 出力関係が非常に非線形であるそれら用には最適変換ではないかも知れない。

【0033】

ここで説明する方法の好ましい実施例では、その組み合わせが "入力フィーチャー (input features)" としても知られる、入力は初期には変換されない。もし次の入力データ集合が、モデル化される必要のある出力に関する十分な情報を現さないならば、上記で説明されたそれらの様なデータ変換が行われてもよい。この戦略を使う主な理由は、変換の形式内に追加的ジオメトリーを課すよりも、可能な所ではどこでも実際のデータを使用する

50

ことである。この追加的ジオメトリーが取る形式は未知であるかも知れない。加えて、データ変換過程を避けることは該変換過程の計算的オーバーヘッドを避け、かくして、特に非常に大きなデータ集合用の計算効率を改善する。

【 0 0 3 4 】

実際のデータが好ましくは変換なしで使用されるのがよいとは云っても、他の入力よりも情報豊富な入力、又はフィーチャーを同定し、選択することにより次元数 (dimensionality) はなお減じられてもよい。これは、入力数が非常に多い時は特に望ましく、最終モデルに起こり得るフィーチャーを全て使用することは非実用的である。データ集合の " 次元 (dimension) " は入力の全部の数として規定されてもよい。実験型モデルを開発する前に、好ましくは、当面のモデリングタスク用に最も情報豊富なフィーチャーを同定されるのがよい。入力数を減じる、又は該問題の次元数を減じる 1 つの技術は、少しの情報コンテンツしか持たない入力を除くことである。これは入力と、対応する出力と、の相関 (correlation) を調べることに依りなされてもよい。しかしながら、好ましくは、次元数削減は、下記で論じる様に、情報豊富と決定されたフィーチャー組み合わせで各入力の発生頻度 (each input's frequency of occurrence) を調べることにより行われるのがよい。それで、より少ない発生頻度の入力 (less-frequently-occurring inputs) はモデル発生過程から排除されてもよい。

10

【 0 0 3 5 】

時間変化する又は動的なシステム用では、追加的複雑さが、与えられた何れかの時の出力が、より早期の時の入力と出力との双方にも左右される事実から生ずる。この様なシステムでは、該データ集合の正しい表現が非常に重要である。もし特定時刻の測定出力に対応する入力とその時だけ測定されるならば、該時間遅れ (time lags) (すなわち、入力発生と該結果としての出力発生との間の時間間隔) 内に含まれる情報は失われる。この問題を緩和するために、入力の拡張された集合から成るデータ表 (data table) が作られるが、そこでは該入力の拡張された集合は入力の現在の集合のみならず多数の前の時刻 (at multiple prior times) の入、出力からも成っている。この新データ表は次いで選択された時刻範囲に亘り (spanning a selected time horizon) 情報豊富な入力組み合わせ用に解析され得る。

20

【 0 0 3 6 】

拡張データ表の創生での重要な事項は時間的に如何に遠くまで逆戻って知るかである。多くの場合、これは先験的には知られず、余りに長く早期までの時間間隔 { 時間範囲 (time span) } を含むことにより、該データ表の次元数は非常に大きくなる。この事項を処理するために、多数のより短い時間範囲のデータ表が元のデータ表から作られるが、各データ表は過去での与えられた時間間隔から成る。これらのより新しいデータ表の各々の及ぶ時間間隔は重なったり、隣接したり又は分離していてもよい。これらのより小さいデータ表の各々からの最も情報豊富な入力が次いで集められ、該小さなデータ表からの選択された入、出力を含むハイブリッドデータ表を作るよう組み合わせられる。この最後のハイブリッド表は、該時間間隔間の起こり得る相互作用が今や含まれるので、次いでデータモデル化過程への入力として使用出来る。

30

【 0 0 3 7 】

例えば、もし住宅販売レート (home sales rate) が商品製材価格 (commodity lumber prices) に影響するが、約 2 ヶ月の推定時間遅れがあるのでないか、を調査したいならば、この時間遅れを発見するために本発明用には該データ表は入力が出力に 2 ヶ月先行する対応 (matched) した入、出力を要する。これは、実際の時間遅れがどれだけかを発見するために種々の入力が 1 つの出力に対し異なる遅れを有する 1 つ以上のデータ表 (すなわち、列は入、出力、行は連続した時間) を形成することにより行われ得る。特に、1 つの出力は X 日の製材価格であってもよい。入力が X 日、X - 1 日、X - 2 日 から X - 120 日までの住宅販売レートであるのみならず、X - 1、X - 2 から X - 120 までからの出力でもある。高い情報コンテンツを持つ最も早期の入力が失われないことを保証するために、入力と対応する出力との間の推定時間遅れ (suspected time lag) より

40

50

長い時間間隔が選択される。次いで次の表の行はY日（例えば、 $X + 1$ 又は幾らかもっと後れた日）の製材価格に等しい出力を有し、入力はY、 $Y - 1$ 、 $Y - 2$ 、...、 $Y - 120$ の住宅販売レートであるのみならず $Y - 1$ 、 $Y - 2$...から $Y - 120$ 日までからの出力でもある。次いで該システムは該出力に影響する入力の組み合わせを同定することにより適当な時間遅れを同定する。2. データ量子化とフィーチャー部分空間内のセル境界一旦適当なデータ表現が確立されると、サンプル点を特徴付けるため使用される各入力では量子化（quantization）過程が行われる。入力値の範囲を部分範囲に分ける、すなわち、当該技術で”ビンニング（binning）”として公知の、ビン（bins）に分けるために2つの量子化方法が使われるが、該ビンニングは与えられたフィーチャー部分空間の各入力で行われるが、そこでは各入力は該部分空間の次元に対応し、それはセルの領域に分けられる与えられたフィーチャー部分空間となる。

10

【0038】

最も簡単な量子化法は固定サイズの部分範囲、すなわちピン幅（時には、”固定ビンニング（fixed binning）”として知られる）に基づくが、そこでは各入力に付随する値の全体範囲が等間隔又は等サイズの部分範囲又はビンに分けられる。

【0039】

もう1つの量子化、それは”統計的量子化（statistical quantization）”と呼ばれてもよく、図2Aで最も良く見られ、ここでは”適合的量子化（adaptive quantization）”と呼ぶが、は値の該範囲を不等サイズの部分範囲に分けることに基づく。もしデータがデータピン210により示す様に均一に分布されていれば、該ピンサイズは大体等しい。しかしながら、該データ分布がクラスター（clustered）されるならば、該ピンサイズは、ピン220により示される様に、各ピンがデータ点の殆ど等しい数を含むように適合的に調整される。図2Bに見られる様に、各部分範囲、又はピンのサイズは、入力範囲を等しい百分位数（percentile）の部分範囲に分け、それらの百分位数を該ピン240を作るフィーチャー値の範囲上に射影（projecting）することにより、各入力の累積確率分布（cumulative probability distribution）230（又はヒストグラム）に関係付けられてもよい。

20

【0040】

この方法で、各入力上のグローバル情報がその入力上で該データを適合的に量子化するため使われる。この方法では、各入力は別々に量子化され、すなわち、量子化は入力毎ベースで行われる。該部分範囲又はピンのサイズ（幅）は与えられた入力内で一般に不均一で、その入力の累積確率分布の形を反映していることを注意すべきである。該部分範囲のサイズは入力から入力へと変わってもよい。適合的量子化（適合的ビンニング）は情報を含まない空の入力の部分範囲を有する確率を減らす、それはさもないと最終モデル内の情報ギャップとなる。

30

【0041】

与えられた入力に対する該部分範囲、又はピンのサイズは部分空間から部分空間へと変わってもよい。すなわち、或る入力は、それらが高い次元の部分空間で現れる時より低い次元の部分空間で現れる時の方がより精細な解像度のビンニングを有してもよい。これは或る全体のセルの解像度（セル当たりの点の数）は、データの意味のある量がセル内で一緒にグループ化又はビン化（binned）されるように、望まれる事実のためである。セル数は次元数に指数関数的に比例するので、より高い次元のフィーチャー部分空間は、セル当たりの望ましい平均の点の数を保持するように、個別入力用により粗いビンニングを使用する。データ量子化がモデル化の方法のローバストさ用に顕著な意味を有するのは該データの残りの外れ値の点の偏差の大きさが該量子化（ビンニング）過程に抑制されるからである。例えば、もし入力値が最高部分範囲（ビン）内の上限を越えるなら、それはその値に無関係にその部分範囲（ビン）内に量子化（ビン化）される。

40

【0042】

ここで使用される”フィーチャー部分空間”は1つ以上の入力の組み合わせと規定される。フィーチャー部分空間の画像的表現が創られてもよく、それも又簡単に”部分空間”と

50

してここでは呼ばれる。該部分空間は好ましくは複数の " セル " に分けられるのがよく、該セルは該フィーチャー部分空間を含む入力の部分範囲の組み合わせにより規定される。好ましい実施例では、データ量子化は更に、(前の説明の固定的か又は適合的か何れかの方法を使用して) 入力当たりの部分範囲 (ピン) の数を規定するか、又は、代わりに、該フィーチャー内のセル当たりデータ点の平均数を規定するか、何れかで指定される。これは適合的量子化法の多次元的拡張と見られる。

【 0 0 4 3 】

図 3 A、3 Bそして3 Cを参照すると、固定サイズのビンニングがそれぞれ 1 , 2そして3次元フィーチャー部分空間で示される。該データ集合は各々が4つの入力、又はフィーチャーを有する4つのデータ点、DP1 - DP4から成る。該データ集合は全ての3つの図で同じである。該データ点はどのフィーチャー (又はフィーチャー組み合わせ) が選択されるかにより特定のセルに分類される。図 3 Aでは、もし該1次元部分空間が第3の入力 (左端のビットに対応する第1入力を用いて0010と呼ばれる) を表せば、DP1とDP4はセルC1に分類され (DP1 = . 5、DP4 = . 3)、DP2とDP3はセルC2に分類される (DP2 = 1 . 2、DP3 = 1 . 7)。もし、しかしながら、該1次元部分空間が第2入力 (0100) であると取られるなら、DP2とDP4はC1に分類され (DP2 = . 7、DP4 = . 4)、そしてDP1とDP3はC2に分類される (DP1 = 1 . 5、DP3 = 1 . 9)。

【 0 0 4 4 】

図 3 Bでは、もし該部分空間が第1と第2入力 (1100) により指定されれば、DP1はセルC2に分類される { DP1 = (. 5、1 . 5) } が、なお該第1と第3入力 (1010) により発生される部分空間ではセルC1に分類される。図 3 Cでは、DP1は第1、第3そして第4入力 (1011) で規定される部分空間ではセルC1に分類され、第1、第2そして第4入力 (1101) で規定される部分空間ではセルC2に分類される。

【 0 0 4 5 】

該入力に基づく該システムの出力の予測で或る精度を有するフィーチャー組み合わせを同定することが望ましい。特定の入力組み合わせ、又はフィーチャー組み合わせは多くのユニークな部分空間を規定することが上記例から分かる。有限数の入力シーケンスを仮定すれば、の部分空間の数は勿論有限であるが、該数は入力数と共に極めて急速に成長する。

【 0 0 4 6 】

フィーチャー選択のタスクは入力 - 入力の相互作用の可能性により複雑化する。このような相互作用が存在すれば、個別には情報貧弱な入力が高い情報エントロピーを有する入力の組み合わせを作る相補的な仕方では組み合わせられ得る。かくして、入力 - 入力相互作用の可能性を無視するどんなフィーチャー選択方法もモデル化過程から有用な入力を排除する可能性があり得る。この制限を避けるために、好ましい方法は、入力 - 入力関係を本質的に含み、該データ内にあるかも知れぬ何等かの非線形性を非常に自然に処理する、情報理論ベースのフィーチャー部分空間を選択する取り組みを使用する。

【 0 0 4 7 】

加えて、該方法は利用可能な部分空間のエグゾースチブ (exhaustive) な探索を含むが、それが好ましくは情報エントロピーのメジャーを適応度関数として使う遺伝的発展型アルゴリズム (genetic evolutionary algorithm) を含むのがよい。

3 . 遺伝的発展と情報エントロピーを使用するフィーチャー部分空間選択

ここで説明する方法は好ましくは " 遺伝的アルゴリズム " として公知の比較的最近のアルゴリズム的取り組みを使用するのがよい。ジョンエイチ・ホランド (John H. Holland) { 1975年発行、アナーバー、ミシガン大学プレス (Ann Arbor:the University of Michigan Press)、" 天然及び人工的システムでの適合 (Adaptation in Natural and Artificial Systems) " で } により定式化され、又デー・イー・ゴールドバーグ (D. E. Goldberg) { 1989年発行、アディソン - ウエズレーパブリッシングカンパニー (Addison-Wesley Publishing Company)、" 探索、最適化及び機械学習に於ける遺伝的アルゴリズム (Genetic Algorithms in Search, Optimization and Machine Learning) " で } 及びエ

10

20

30

40

50

ム・ミッチェル (M. Mitchell) { 1997 年発行、エムアイテーパーレス (M.I.T. Press)、" 遺伝的アルゴリズム入門 (An Introduction to Genetic Algorithms) " で } により説明された様に、該取り組みは最適化問題を解く強力で、一般的な方法である。遺伝的アルゴリズムの取り組みは次の様である。

【 0048 】

(a) 問題の解空間 (solution space) を N ビット記号列 (N-bit strings) の母集団 (population) としてエンコードする。ポピュラーなエンコード用フレームワークは 2 進記号列 (binary strings) に基づく。該ビット記号列の集まりは " 遺伝子プール (gene pool) " と呼ばれ、個別ビット記号列は " 遺伝子 (gene) " と呼ばれる。

【 0049 】

(b) 目前の問題に対する何等かのビット記号列の適応度 (fitness) を測定する適応度関数 (fitness function) を規定する。換言すれば、該適応度関数は何等かの起こり得る解の良さ (goodness) (又は精度) を測定する。

【 0050 】

(c) ビット記号列のランダムな遺伝子プールで最初にスタートする。それを通してより " 適した (fit) " ビット記号列が " より適した (fitter) " 子供 (offspring) の新しいプールを作るために優先的にメートする、選択的再組み合わせ (selective recombination) 及び突然変異 (mutation) の様な、遺伝子から得られたアイデアを使用することにより、より適したビット記号列の次の世代が発展出来る。 " 適応度 (Fitness) " は情報エントロピーのメジャーにより決定される。突然変異の役割は起こり得る解の探索空間を拡張することであり、該解は改善された度合のローバストさ (robustness) を創る。

【 0051 】

(d) 上記進め方に従う数世代の発展の後、より適したビット記号列のプールとなる。最適解はこのプール内の " 最適 (fittest) " ビット記号列として選択される。

【 0052 】

これらの側面の各々を下記で更に詳細に論じる。

a . N ビット記号列の母集団としての解のエンコーディング (Encoding solution as a population of N-bit strings)

最適問題を解くために遺伝的アルゴリズムを使う最初の過程は、ビット記号列として表される解となる方法で該問題を表すことである。簡単な例は 4 入力と 1 出力を有するデータベースである。入力の種々の組み合わせが 4 ビット 2 進記号列により表される。該ビット記号列 1 1 1 1 は、全ての入力为该組み合わせ内に含まれる入力組み合わせ、又はフィーチャー部分空間を表す。最左ビットを入力 A、第 2 の最左ビットを B、第 3 の最左ビットを入力 C そして最右ビットを入力 D と呼ぶ。もしビットが値 1 に換わるなら、それは対応フィーチャーが該組み合わせ内に含まれるべきことを意味する。逆に、もしビットが値 0 に換わるなら、それは対応フィーチャーが該組み合わせ内で排除されるべきことを意味する。

【 0053 】

同様に、該ビット記号列 1 0 0 0 は唯フィーチャー A が含まれ、全ての他の入力が排除される入力組み合わせを表す。この方法で、16 の全可能性からのあらゆる起こり得る入力組み合わせは 4 ビット 2 進記号列により表される。一般に、もしモデル化されるデータベースに N 入力があるなら、全ての起こり得る入力組み合わせは N ビット 2 進記号列を使用して表される。4 次元のフィーチャー部分空間を表すサンプルの 2 進ビット記号列は図 4 に示される。図 4 の該ビット記号列は D ビットを有し、その 4 つだけが " 1 " のビットである。該 " 1 " のビットは 4 つのフィーチャー F1, F4, Fi、そして FD と対応する。該変数 i と D は一般化された場合を表すために使用される。更に進んだ例が図 3 A で示されるが、そこでは 4 入力システムを表し、1 つの " 1 " ビットを有する、4 ビット記号列が 1 次元フィーチャー部分空間に対しコード化する。2 つの " 1 " ビットが図 3 B に見られる 2 次元部分空間に対しコード化し、3 つの " 1 " ビットが図 3 C で見られる 3 次元部分空間に対しコード化する。

10

20

30

40

50

b. ビット記号列の適応度を測定するための適応度関数の規定

最適化問題への解として最適ビット記号列を発展させるために、発展過程をドライブするため使用される定量評価 (metric) を規定することが必要である。この定量評価は遺伝的アルゴリズムでは適応度関数と呼ばれる。それは与えられたビット記号列が如何に良く目前の問題を解くかのメジャー (measure) である。適当な適応度関数を規定することは該ビット記号列がより良い解へ発展することを保証する重要過程 (critical step) である。

【 0 0 5 4 】

上記例では、各 4 ビット 2 進記号列は入力 of 起こり得る組み合わせをエンコードする。入力フィーチャー部分空間は、対応するビット記号列内でオンに換わる入力フィーチャーを使用することにより作られ得る。データベース内のデータはこのフィーチャー部分空間内へ射影され得る。該適応度関数は、該入力フィーチャー部分空間上で出力状態の分布を調べることににより情報豊富さのメジャーを提供する。もし該出力状態がこの部分空間上で非常にクラスターされてそして分離されていれば、該対応する入力フィーチャー組み合わせは異なる出力状態を分離することでよい仕事をしているので該適応度関数は高い値となる。逆に、もし全ての出力状態が該部分空間上にランダムに分布されているならば、該対応する入力フィーチャー組み合わせは該異なる出力状態を分離することで貧弱な仕事をしているので該適応度関数は低い値となる。代わりに、該適応度関数は、該部分空間内の個別セルの情報豊富さを調べ、次いで該セルの加重平均を形成することにより該部分空間の情報豊富さのメジャーを提供してもよい。

【 0 0 5 5 】

好ましくは、出力状態クラスタリングのグローバルなメジャーは最良のビット記号列の発展をドライブする該適応度関数として使用される。このメジャーは好ましくはクラスタリングを規定する強力な方法であるエントロピー関数に基づくのがよい。適応度関数のこのエントロピー的規定を用いて、該出力を最も良くクラスターし分離する入力組み合わせを表すビット記号列が該発展型過程から出現する。代替りの適応度関数は、出力状態確率の標準偏差か分散か、又は少なくとも 1 つの出力確率が他の出力確率より顕著に大きい部分空間内のセル数を表す値かを含む。出力状態の集中を測定する他の同様な発見的方法 (heuristics)、又はアドホック (ad hoc) な規則は発展型過程内で容易に交換される。

c. 発展型過程の詳細

1. N ビット 2 進記号列のランダムなプールの創生

図 5 A を参照すると、該発展型過程 5 0 0 は過程 5 1 0 で始まり、そこでは N ビットの 2 進記号列のランダムなプールが創られる。これらの初期 2 進記号列は、それらがともかく最適であると言う先験的理由がないので一般的にそれらの適応度関数用には非常に低い値しか持たない入力フィーチャー組み合わせをエンコードする。この初期プールは該発展型過程を開始するため使われる。

【 0 0 5 6 】

2. 適応度の計算

該プール内の各 2 進記号列の適応度は過程 (b) で説明した方法を使用して計算される。該データは過程 5 2 0 で示すようにバランスを取られる。各 2 進記号列用にフィーチャー部分空間が発生され、データベース内のデータが対応する部分空間内へ射影される。該部分空間は過程 5 3 0 で行われた選択に従って、等間隔のピンニング 5 3 2 又は適合的に隔てられたピンニング 5 3 4 の選択に依りピンに分けられる。考慮下の特定の遺伝子が過程 5 4 0 で選択され、そしてピンの数は過程 5 5 0 で、好ましくはユーザー入力により、ピンの固定数 5 5 2 を指定するか又はセル当たりサンプルの平均数 5 5 4 を指定することにより決定される。該ピン配置は次いで過程 5 6 0 に示す様に、決定される。次いで対応 2 進記号列の適応度を表す出力状態のクラスタリングと分離の程度を計算するためにエントロピー関数又は他の規則が使用される。これは、データ点が各部分空間内に配置される過程 5 7 0 と、グローバル情報コンテンツが決定される過程 5 8 0 で示される。過程 5 8 5 により示される様に、次の遺伝子シーケンスは過程 5 4 0 の開始で動作する。

【 0 0 5 7 】

3. 適応度の加重ルーレットホイール (weighted roulette wheel) の創生

各2進記号列の適応度が計算された後、該適応度の加重ルーレットホイール592が図5Cに示す様に創られる。これは、より高い適応度値 (fitness value) を有する2進記号列がより低い適応度値を有する2進記号列よりも比例してより広いスロット幅に付随される過程と考えられる。これは、該ルーレットホイールが廻されると、より低い適応度の2進記号列よりも、より高い適応度の2進記号列の選択に、より重く加重する。この過程は下記で更に詳細に説明する。

【0058】

4. 新しい親の2進記号列 (new parent binary strings) の選択

ルーレットホイール592は次いで廻され、該ホイールが終わるスロットに対応する2進記号列が選択される。もし元のプールにN個の2進記号列があるなら、該ホイール592はN個の新親記号列を選択するためN回廻される。ここで重要な点はもしそれが高い適応度値を有するなら該同じ2進記号列が1回より多く選ばれ得ることである。逆に、低い適応度関数を有する2進記号列は、それが完全に排除されることはないが、親として決して選択されないことが起こり得る。次いでN個の親が、新しい子の2進記号列発生への先駆者としてN/2個の対に対化される。

【0059】

5. 子記号列を創る親の交叉 (crossover) と突然変異 (mutation)

一旦2つの親が選ばれると、図5Dに示す、交叉オペレーション (crossover operation) 594が行われるべきか否かを決定するために加重コインがフリップされる。もしこれが交叉オペレーションとなるなら、クロッシングサイトがビット位置1と該記号列内の最後のビット位置の次にあるの最後の起こり得るクロッシングサイトとの間でランダムに選択される。該クロッシングサイトは各親を右側と左側に分割する。図5Dに示す様に、各親の左側を他の親に右側と連結することにより2つの子記号列が創られるが、そこでは該親遺伝子10001と00011は左半分100と000、そして右半分01と11に分割され、次いで10011と00011を形成するよう組み合わせられる。最後に、該2つの子記号列が創られた後、該子記号列プールの多様性を増やすために該子記号列の小数の個別ビットがランダムに逆にされる (突然変異される)。これは与えられたビットが逆にされる確率に換算して指定出来る。逆転の確率は望ましいビット突然変異の数と該記号列内ビット数に基づいて尺度合わせされる。すなわち、もし記号列当たり平均5つの突然変異が望まれるならば、与えられたビット変更の確率は100ビット記号列用に0.05に、そして50ビット記号列用に0.1等に設定される。

【0060】

6. 発展型過程の継続

過程590に示す様に、上記過程2-5は、各創られた子記号列プールを次世代用の新しい親プールとして使用して、数回 (又は数世代) 繰り返される。該子記号列プールが発展すると、それらの対応適応度は平均で改善すべきであるが、それは各世代で、新しい子記号列を創るために、より適した記号列が優先的にメートされるからである。

【0061】

該発展型過程は、予め決められた数の世代の後か、又は最高適応度の記号列か又は平均プール適応度か何れかが最早変化しない時か、何れかで停止出来る。

【0062】

最適化問題を解くための遺伝的アルゴリズムの使用で、解かれる必要にある2つの重要な項目がある。第1の項目はエンコーディングスキームである。該問題がビット記号列としてエンコードされ得る解の役に立つか? 第2の項目は該適応度関数の選出である。該発展型過程は該適応度関数により統制される (すなわち、導かれる) ので、その解の質は間近な目標への適応度関数のマッチングに密接に依存している。

【0063】

ここに説明した好ましい方法では、第1の項目は、図4で図解され、各ビットがデータ集合のNの入力の1つと対応する、Nビット2進フイーチャービット記号列を含む遺伝子を

10

20

30

40

50

規定することにより解決される。該Nビット2進フイーチャービット記号列の各ビットは対応入力を参照し、もし該対応入力があるフイーチャー部分空間内であれば該値1を、もし該対応入力があるフイーチャー部分空間内に無ければ該値0を有する。

【0064】

該好ましい方法では、第2項目はフイーチャー部分空間のグローバルエントロピーを計算する情報エントロピーメジャー (informational entropy measures) を使用することにより解決される。該フイーチャー部分空間のグローバルエントロピーは、それから最適モデルが発展させられ得る最適フイーチャー組み合わせのプールの発展をドライブする適応度関数として使用される。該グローバルエントロピーは、フイーチャー部分空間内のセルのローカルエントロピーを最初に決定し、そして該ローカルエントロピーの加重和として全体のフイーチャー部分空間のグローバルエントロピーを計算することにより計算される。代わりに、部分空間のグローバルエントロピーは、該全体の部分空間の間で、与えられる出力用の点の分布を調べ、そして次いで全ての状態に亘り特定状態向けエントロピーの加重平均を形成することにより決定されてもよい。フイーチャー部分空間プールを保持する能力は、そのどちらも最終モデルのローバストさに寄与する該解空間内の冗長度と多様性の双方を提供する。

ローカルセルエントロピーとグローバル部分空間エントロピーの決定

好ましい方法の側面に依れば、情報コンテンツのレベルが測定される。特に、セル又は部分空間の情報コンテンツのレベルはデータ分布の均一性のメジャーである。すなわち、データが均一である程、システムのモデル化の目的にそれが持つ予測価値は大きくなり、従って、情報コンテンツのレベルは高くなる。該均一性は多数の代替的方法で測定されてもよい。1つのこの様な方法はクラスタリングパラメーター (clustering parameter) を使用する。用語クラスタリングパラメーターはローカルセルエントロピー、考慮下の特定部分空間上で計算された特定出力のエントロピー、又はここで論じられる発見的方法、又は他の同様な方法を指す。

【0065】

図6を参照すると、個別セルの情報コンテンツは方法600により示されたカテゴリー的出力システム及び方法602による連続する定量的モデル用に決定される。好ましい実施例では、前に論じたニシ (Nishi) の情報エントロピー規定が、該情報コンテンツを表すローカル及びグローバル両エントロピー加重を数学的に規定するため使用される。本発明の実験型モデリング用には、ニシにより拡張された、シャノンのエントロピーの概念が、該エントロピーのメジャー (measure) が計算されるデータ集合用の適当なメジャーであることが見出されて来た。ニシの式が出力状態に対応する確率の集合に適用される。等しい出力確率を有するセル (各出力が等しく似ている) は少しの情報コンテンツしか有しない。かくして、高い情報コンテンツを有するデータ集合は他より高い、幾らかの確率を有する。より大きな確率的変動 (greater probabilistic variations) は出力状態の不均衡 (imbalance in the output states) を反映し、従って該データ集合の高い情報豊富さの指標を与える。

【0066】

好ましい方法では、一般的なエントロピー加重項 (general entropic weighting term) W が規定され、 $W = 1 - E$ の形式を有する。該エントロピー加重項 W はニシの情報エントロピー関数 E の補数 (complement) であり、完全に不均一な分布用に値1を有し、完全に均一な分布用に値0を有する。

【0067】

図6の方法600を再び参照すると、情報レベルはローカルエントロピー加重項 (local entropic weighting term) を計算することにより決定される。例えば、部分空間内の与えられたセル用に適当なものは次の仕方で規定され得るが、すなわち最初に、過程610で、 n_c エントリーを有するデータ集合が創られ、ここで n_c は出力状態の数である。各エントリーは下記で与えられるセル i 用の特定状態向けローカル確率 $p_{c|i}$ に対応しており

、

【 0 0 6 8 】

【 数 5 】

$$p_{ci} = n_{ci} / \sum_{k=1}^n n_{ki},$$

【 0 0 6 9 】

ここで n_{ci} は c の出力状態を有するセル i 内の点の数であり、該和はセル i 内の全ての出力状態 k に亘り延び、かくしてセル i 内の全ての点を含む。与えられセル i 用に、値 p_{ci} のシーケンスは種々の出力状態 c にある確率を表す。過程 6 2 0 で該セルの情報コンテンツは決定される。好ましくは、ニシの情報エントロピー規定が部分空間 S 内の与えられたセル i 用のローカルエントロピー項 E を規定するため使用されるのがよく、

10

【 0 0 7 0 】

【 数 6 】

$$E_i^S = (\sum_{k=1}^{n_c} f_{k|i}^S \ln f_{k|i}^S) / \ln(1/n),$$

【 0 0 7 1 】

ここで和の変数 k は出力状態、 n_c は出力状態（又は“カテゴリー”）の総数を表し、そして

20

【 0 0 7 2 】

【 数 7 】

$$f_{k|i}^S = p_{ci}^S / \sum_{k=1}^{n_c} p_{ki}^S.$$

【 0 0 7 3 】

である。

【 0 0 7 4 】

勿論、全ての k に亘る全ての $p_{k|i}$ の和は 1 に等しいが、明確化のため上記に含まれる。

30

【 0 0 7 5 】

最後に、又過程 6 2 0 で、該ローカルエントロピー加重係数は

$$W_i^{L^S} = 1 - E_i^S$$

であり、ここで上書き L^S は W が部分空間 S 内でセル用のローカルエントロピー関数であることを呼称する。高い情報コンテンツを有するセルは高いローカルエントロピー加重を有する。すなわち、それらは $W_i^{L^S}$ の高い値を有する。

【 0 0 7 6 】

代わりに、該情報コンテンツは、該出力確率値の分散又は標準偏差を決定することによるか、又は何等かの 1 つの出力が予め規定されたしきい値を上回る付随確率を有するかどうかを決定することによる様な、均一性のもう 1 つのメジャーにより測定されてもよい。例えば、セルの確率分布に基づきセルに値を割り当ててもよい。特に、予め決められた値より大きい何等かの出力状態確率を有するセルは 1 の値を割り当てられ、該出力状態確率のどれも予め決められた値より大きくないどのセルも値 0 を割り当てられる。該予め決められた値は該フィーチャー部分空間（モデル、フレームワーク、スーパーフレームワーク等）の結果に基づき実験的に選ばれた定数である。該定数は又出力状態の数に基づいてもよい。例えば、何れかの出力状態が平均より大きい発生の尤度（greater-than-average likelihood of occurring）を有するセルの数を数えたいと願ってもよい。それで、 n の出力状態システムについて、 $1/n$ より大きい何等か 1 つの出力状態確率を有するどんなセルも 1 の値を与えられるか、又は k/n より大きければ、或る定数 k が与えられる。他のセル

40

50

はゼロの値を与えられる。

【 0 0 7 7 】

代わりに、セルに与えられる加重は与えられた確率を越える出力状態の数に基づいて増加出来る。例えば、4出力状態システムでは、0.25より大きい発生確率を有する2つの出力状態を有するセルは2の加重を与えられる。更に進んだ代替えとして、セルの又はグローバルな加重は出力状態の分散に基づくことが出来る。他の同様な発見的方法が考慮下のセルの情報コンテンツを決定するため使用されてもよい。

【 0 0 7 8 】

モデル化されつつある過程の出力が連続的な場合、ローカルエントロピーは方法602に示す様に計算される。過程630で、該セルに存在する出力値の全てを含むデータ集合が創られる。該セルの情報コンテンツは過程640で計算される。出力に特定の確率を処理する時、高い情報コンテンツを有するデータ集合は他より高い或る確率を有することが思い出される。出力値を直接処理する時、しかしながら、過程630 - 670でその場合である様に、情報豊富な集合はより均一なデータ値を有するそれらである。すなわち、高い情報集合は出力値ではより少ない変動を有する。かくして、もし情報コンテンツが該ニシのエントロピー計算を使用して決定されれば、該補数的値 $1 - E$ を形成する必要はない。この場合の加重係数は簡単にニシのエントロピー E に等しい。

【 0 0 7 9 】

加えて、過程650と660で示す様に、低エントロピーセルにゼロを設定するようにしきい値限定を適用することが望ましい。これはグローバルな計算が行われる時意味のない情報コンテンツを有するセルの情報コンテンツを累積することに付随する誤った影響を制限する助けになる。ローカルなセルのエントロピーの計算は過程670に示す様に完了する。

【 0 0 8 0 】

代わりに、連続的出力システムを取り扱う時、該出力を複数のカテゴリーに量子化し、各量子化レベルでの確率を有するデータ集合を規定するために、過程610で示す上記方法の過程を使用することが可能である。残りの過程620も、上記説明の様にエントロピー加重を計算することによって、該情報コンテンツを決定するため行われる。

ローカルエントロピーの加重和としてのグローバルエントロピーの計算

図7を参照すると、部分空間 S 用のグローバルエントロピー W^{gs} は次いで、その部分空間内の全セルに亘りローカルセルエントロピー W^{ls} のセル母集団加重和 (cell-population-weighted sum) として計算される。

【 0 0 8 1 】

【 数 8 】

$$W^{gs} = \sum_{i=1}^n n_i^s W_i^{ls} / \sum_{i=1}^n n_i^s,$$

【 0 0 8 2 】

ここで n は部分空間 S 内のセル数を表し、 n_i^s は部分空間 S 内のセル i 内のカウント (データ点) 数を表す。実際は、これは、それがその部分空間内のセルのプューリテイ (purity) の全体的メザーを記述するので、グローバルエントロピーの有用なメザーであることになった。図8はローカルとグローバルの情報コンテンツの計算を図解する。図9はローカルとグローバルのエントロピーパラメーターの例を示す。高い情報コンテンツを有する部分空間は W^{gs} の高い値を有する。出力状態依存のグローバルエントロピーを計算する代替え的方法

規定された基本的統計量は、該出力が部分空間 S 内の状態 c 内にあるとした場合にセル i 内にある確率を表す確率 $p_{i|c}$ である。

【 0 0 8 3 】

【 数 9 】

10

20

30

40

50

$$p_{ic}^s = n_{ci} / \sum_{j=1}^n n_{cj},$$

【 0 0 8 4 】

ここで n_{ci} は出力状態 c を有するセル i 内の点の数であり、該和は部分空間 S 内の全てのセル j に亘って伸展する。

【 0 0 8 5 】

該ニシの情報エントロピー規定が部分空間 S 内の与えられた出力状態 c についてグローバルエントロピー項 W_c^{gs} を規定するため使用出来る。最初に、与えられた状態 c 用のニシ

10

のエントロピーが計算される：

【 0 0 8 6 】

【 数 1 0 】

$$E_c^s = (\sum_{i=1}^n f_{ic}^s \ln f_{ic}^s) / \ln(1/n)$$

【 0 0 8 7 】

ここで n はセル数であり、

【 0 0 8 8 】

20

【 数 1 1 】

$$f_{ic}^s = p_{ic}^s / \sum_{j=1}^n p_{jc}^s.$$

【 0 0 8 9 】

である。

【 0 0 9 0 】

再び、状態に特定の確率 (state-specific probabilities) の全てのセルに亘る和である、分母は 1 に等しいが、一貫性と明確化のために上記表現に含まれる。 E_c^s はかくして該部分空間 S 上の確率 p_{ic}^s の分布のグローバルな均一性を表す。最後に、該グローバルエントロピー項 W_c^{gs} は下記で規定され

30

$$W_c^{gs} = 1 - E_c^s.$$

それは部分空間 S 内でのカテゴリー c 用のグローバルな出力に特定のエンロピー加重項である。これは、それが全体の部分空間を通しての点の分布 (出力 c に対応する) のクラスタリングを表す意味でグローバルなメザーである。高い情報コンテンツを有する部分空間は高い値の W_c^{gs} を有する。

グローバルエントロピー加重係数の代替的規定用のカテゴリーから独立した一般化

全カテゴリーに亘り加算することにより、代替的グローバルエントロピー加重係数はカ

40

テゴリーから独立したグローバルエントロピー加重係数として規定され

【 0 0 9 1 】

【 数 1 2 】

$$\bar{E}^s = (\sum_{c=1}^n \sum_{i=1}^n f_{ic}^s \ln f_{ic}^s) / \ln(1/n')$$

【 0 0 9 2 】

ここで n' は $= n_c \cdot n$ で、それは出力状態数とセル数の積であり、ここでは

【 0 0 9 3 】

50

【数 1 3】

$$f_{ic}^S = p_{ic}^S / (\sum_{c=1}^{n_c} \sum_{i=1}^n p_{ic}^S).$$

【0 0 9 4】

である。勿論、上記式の分母は

【0 0 9 5】

【数 1 4】

$$\sum_{c=1}^{n_c} \sum_{i=1}^n p_{ic}^S = n_c$$

10

【0 0 9 6】

と簡単化され、それはニシの式で使用される確率が適切に正規化されることを示す。この代替えの規定は出力状態数が多く、そして計算効率が望まれる状況で有用と信じられる。

【0 0 9 7】

上記議論で、該システムの出力値が離散的 (discrete)、又は "カテゴリー的 (categorical)" であることが仮定されている。同じ方法は、エントロピー計算の前に最初出力値を離散的状态又はカテゴリーに人工的に量子化することにより、例え該出力値が連続的であつても、ローカル及びグローバルエントロピーを計算するため使用される。

20

【0 0 9 8】

トレーニングのデータ集合の出力状態の母集団の分布は該モデルの究極的有效性 (ultimate validity) に付随されることは述べる価値がある。上記解析で、該データ集合はバランスされていると仮定されてもいるが、しかしながら、この様なことは常にはその場合ではない。2つの出力状態、AとBとがある問題を考える。もし該トレーニングデータ集合が状態Aを表すデータ項目から主として成るならば、該母集団の統計はアンバランスとなり、ことによると偏倚されたモデルの創生となる。インバランスの理由は、データコレクター (data collector) の部分での偏倚か、又は該データ集合の親母集団特性にある真性のインバランスか何れかである。

30

【0 0 9 9】

該データコレクターの部分での偏倚の場合、セル内の母集団統計がデータ項目の絶対数より寧ろ該セル内に存在する与えられた出力状態のデータ項目の部分参照するように簡単な正規化が行われ得る。この正規化は多くの実験データ集合で成功裡に使われて来た。第2の場合では、該インバランスは "真実 (real)" であるので、正規化は適當ではないかも知れない。

【0 1 0 0】

データ正規化の例は次の様である。

【0 1 0 1】

2つの出力状態AとBがある100項目を有するデータ集合を考える。状態Aに対応する75項目と状態Bに対応する25項目とがあると仮定する。状態Aに対応する5項目と状態Bに対応する5項目を有する全部で10項目がある部分空間内のセルを考える。絶対項では、我々は各エントリが特定の状態用のカウントを参照する { 5, 5 } に対応する "カウントデータ集合" を有するので、これはインピュアセル (impure cell) である。しかしながら、該データは次の様にその状態用の全体のカウントに対して各カウントを正規化することによりバランスさせられてもよい。

40

【0 1 0 2】

【表 1】

状 態	カウント	全体の分数
A	5	$5 / 75 = 1 / 15$
B	5	$5 / 25 = 1 / 5$

【 0 1 0 3 】

該表からの該分数的カウントは次いでエントロピー計算で使用される。

10

【 0 1 0 4 】

データ集合Dは $D = \{ 1 / 15, 1 / 5 \}$ 、 $d_{total} = 1 / 15 + 1 / 5 = 4 / 15$ を伴い、正規化されたデータ集合Fは $F = \{ 1 / 4, 3 / 4 \}$ となる。エントロピーEは次の様に計算される。

【 0 1 0 5 】

$E = \{ 0.25 \ln(0.25) + 0.75 \ln(0.75) \} / \ln(1/2) = 0.811$

変型されたニシのエントロピーWは $1 - E$ 、すなわち $1 - 0.811 = 0.189$ である。図2Cはデータ集合内で与えられた出力状態が支配的な時データの影響をバランスさせる方法を図解するブロック図である。

20

予測指向の適応度関数を用いたモデル発展

一旦入力量子化され、フィーチャー部分空間のプールが遺伝的アルゴリズムにより初めに同定されると、それらの好ましい部分空間の組み合わせを形成することによりモデルが発生される。上記説明の様に、データ又はトレーニングデータ集合と呼ばれるデータの部分集合は、そこから情報が抽出され得る多くのフィーチャー部分空間トポグラフィ (feature subspace topographies) を創るために使用される。高い情報コンテンツを有する部分空間が一旦同定されると、これらの部分空間は、出力予測の目的で該データが内部へ射影される " ルックアップ (look up) " 部分空間として使用される。

【 0 1 0 6 】

特定の部分空間による出力予測は該特定の部分空間内の与えられたセル内の出力状態の分布により決定される。すなわち、各データ点 (又はテストデータ部分空間内の各点) は、図3A-Cに関係して見られる様に、与えられた部分空間内の1つのセル内に分類される。各データ点に付随する出力を予測しようとして、人は、部分空間 (全体のデータ集合、又はトレーニング部分集合) を占めるため使用されるデータの分布を単に見て、予測に到達するためこれを使用する。特定の部分空間による出力予測用に従う簡単な規則は、該出力が状態cにあるとなるべき確率が $p_{c|i}$ により与えられることである。この " ローカル " 確率はフィーチャー部分空間内の与えられたセルを占めるサンプル点の出力分布を単に表している。

30

【 0 1 0 7 】

与えられたモデルは部分空間の組み合わせであり、従って、該モデル内の考慮下の全ての部分空間に関して各点が調べられる。該ローカル確率は本質的に " ベース (base) " 量であり、それは次いでモデル内のローカル及びグローバルの両エントロピーにより加重される。該用語 " ローカルエントロピー " と " グローバルエントロピー " は " エントロピー的係数 " 又は " エントロピー的加重 " としてここでは集散的に引用される。それは、簡単な確率的モデルと比較した時本方法をかなりより精密化するモデル予測を決定するグローバル及びローカルの両方の情報定量評価 (information metrics) の追加である。このエントロピー係数の目的は " 情報豊富 " な部分空間内の " 情報豊富 " なセルを際立たせ (emphasize)、個別的に情報が貧弱か { すなわち、情報豊富さの少ない (less information-rich) }、又は情報貧弱な部分空間内に置かれるか何れかであるセルを軽視 (de-emphasize) することである。

40

50

【0108】

かくして発展型モデル過程をドライブするため使用される各部分空間組み合わせ又はモデル用の適応度関数は、予測のエントロピー的加重和と、該予測と該テストデータ点に付随する実際の出力値との間の付随誤差率 (associated error rate) とである (再び、全体データ集合か又は部分集合かの何れか)。

【0109】

かくして、該方法の1側面に依ると、ローカル及びグローバルエントロピー加重係数は該フィーチャー部分空間の情報コンテンツを特徴付けるために使用される。フィーチャー部分空間セルの寄与をローカル及びグローバルな情報メザーにより加重することにより、該方法は種々の種類のノイズ源を有効に抑制することが出来る。1つのこの様なノイズ源はセル内のローカルノイズである。もしセル内の出力状態の分布が均一であるなら、そのセルは少しの予測情報しか有しない。与えられた出力状態の確率はセル内の出力状態の全分布の性質をほのめかすことは出来るが、それは全体の物語は述べない。全ての他の出力状態の分布は与えられた出力状態の確率内には含まれない。2進出力システムの他の何れでも、1つの出力状態確率内に含まれた情報はかくして不完全である。個別セルに付随するローカルエントロピー項の計算は全体のローカル確率分布を特徴付ける加重係数となる。

【0110】

上記説明の様に、該グローバルエントロピー係数は比較目的に幾つかの異なる方法で計算出来る。部分空間のグローバルエントロピーを規定する好ましい技術はグローバルエントロピーをローカルセルエントロピーのセル母集団加重和 (cell-population-weighted sum) として規定することである。該ローカルエントロピーは部分空間内の各セル用に計算され、この部分空間用の該グローバルエントロピーは次いで全てのセルに亘りセル母集団加重和を行うことにより計算される。これは部分空間について全体のグローバルセル情報エントロピーを測定する (部分空間のセル全部上で)。

【0111】

代わりのグローバルメザーは全体の部分空間上で該セル内の各出力状態の確率分布を調べる。もしこの分布が均一なら、関心のある該部分空間はその出力状態について少しの予測情報しか有さない。この実施例で、部分空間内で各出力状態用に別々のグローバルエントロピー項が計算される。この代わりのグローバルエントロピー項は、各出力状態用に同じである、前に説明したグローバルエントロピー項とは異なる。この代わりのグローバルエントロピーのメザーは、与えられた部分空間が1つの出力状態に関しては "情報豊富" であるが、異なる出力状態に関しては "情報が貧弱" である可能性を受け入れる。

【0112】

本方法はノイズを抑制するためにローカル及びグローバルの両方のベースの加重係数の独立した計算を考慮する。これらの係数は最大の予測精度用にローカル及びグローバル情報の間の最適バランスを得るために個別に調整、又は "ツイーク (tweaked)" される。多くの従来技術のデータモデリングシステムでは、ローカル及びグローバル加重係数の相対的大きさを便利に調整することは難しい。前記の様に、大抵の従来技術の方法は解に到達するために全体のデータ集合上での目的関数 (objective function) の最適化に依存する。

【0113】

もう1つの関連項目は冗長度 (redundancy) のそれである。幾つかの入力フィーチャーは与えられた出力に関する本質的に同じ情報コンテンツを含んでいる。例え2つのフィーチャーが特定の出力状態に関する情報を含まなくても、それらはなお相関しているかも知れない。冗長度は本発明の方法を本質的に制限せず、事実、それは全体の計算コストを増やすけれども、創られるローバストさを該モデルに組み入れる方法として非常に役立ち得る。情報メザーを使用するクラスタリング方法はフィーチャー間の冗長度を同定するために利用可能であり、下記で論じる。

【0114】

ローカル及びグローバルの両方のエントロピー加重係数は分布の "構造" 量 (amount of

10

20

30

40

50

"structure")を測定する。分布がより少ししか均一でない、又は"より多く構造化されて(more structured)"いる程、その対応するエントロピー加重 W はより高い。データ空間の構造のこの側面はローカル及びグローバルの統計の重要性を加重するため使用される。

【0115】

ローカル及びグローバルの両エントロピー項の計算は該方法でのローカル及びグローバルな情報加重係数の別々な制御を考慮する。生ずる自然な問題はローカルさの規定であり、ローカルとはどれ程ローカルなのか？この質問の回答は勿論取り組まれる特定の問題による。好ましい実施例に依れば、該方法は該ピンの解像度を走査することによりローカルさの最良の説明を系統的に探索するが、該解像度は今度は最高の予測精度を提供するために多次元のセルサイズを決定する。特に、情報豊富なフィーチャー部分空間の異なるグループが同定され(エグゾースチブな探索か又はフィーチャー部分空間発展かの何れかにより)、そこでは各グループは部分空間当たり異なる数のセル n を使用する。事実、セル数 n は最小値から最大値までエグゾースチブに探索される。セルの最大数はセル当たりの点の最小平均の意味で指定されるが、それは余りに多くのピンで部分空間の分解能を上げ過ぎることは望ましくないからである。最小数は1より例え小さくてもよい。

【0116】

この点で出力状態の特性をより詳細に考慮することは余談に入る価値がある。本発明の方法では、入力の量子化は多次元部分空間を創るために行われる。分類問題では、該出力変数は離散のカテゴリー又は状態であり、かくして既に量子化されている。定量的モデリングでは、出力変数は連続的である。この様な場合、1つの起こり得る解は該出力状態空間の離散ピンへの人工的な量子化を行うことである。該出力データ空間が量子化された後、上記で説明した離散的モデリングフレームワークがローカル及びグローバルエントロピー係数を測定するために使用され得る。これらのエントロピー係数は下記説明の方法を用いて該出力の連続値の予測に使用され得る。

【0117】

精度に関する重要なメザーは出力状態カテゴリーの数、 n_c の平均全セル母集団統計に対する比 $<n_{pop}>$ である。もし n_c が $<n_{pop}>$ より遙かに大きければ、大抵の出力状態はセル内で空いており、貧弱な統計となり、モデルでの起こり得る劣化となる。これは再び多くのデータを主張し(argues for)、それはデータドライブされるモデルには当然である。コンピュータハードウェア技術の進歩と共に、多量のデータ集合の取得と記憶の能力は急激に増加し、本発明の方法は該データからの情報抽出を可能にする。該方法は、 n_c の値が小さい(1 - 10の桁で)多くの真実の世界の問題で n_c が $<n_{pop}>$ より遙かに大きい時でも驚く程良く作動することが分かった。これは多数の部分空間上での加算統計の協力効果のためかも知れない。

【0118】

抄録すると、フィーチャー部分空間に付随するグローバルエントロピー係数は、遺伝的アルゴリズムを使用して最も情報豊富なフィーチャーのプールを発展させるため使用される適応度関数として使用され得る。このプールの決定は前に説明したデータ量子化条件に依存する。セル当たりサンプル点の平均数が減少すると、該ローカル及びグローバルエントロピー情報メザーは一般に増加する。しかしながら、これは、これらの量子化条件が最終モデルの開発で良く一般化することを必ずしも意味しない。実際に、セル当たりサンプル点の平均数が1より可成り少ない(すなわち、0.1以下)量子化条件下でフィーチャーを発展させることはなお精確なモデルに帰着する。これは主に、該フィーチャープール内の多数の部分空間上での加算統計の協力効果のためである。

システム入力からシステム出力を最も精密に予測するフィーチャーデータ集合の部分集合の決定

図10を参照すると、高い情報エントロピーを有するフィーチャーデータ集合が一旦決定されると、このフィーチャー集合は予測モデルを直接開発するため使用されてもよい。しかしながら、発展型方法(evolutionary method)を使用する該フィーチャー選択過程は

、比較的高い情報エントロピーを有する高次元数データ空間内でそれらのフィーチャーのみを保持することによりいわゆる“次元数の災い (curse of dimensionality)”を緩和する可成りの利点を有する。この関係で、N次元空間内の起こり得る2進フィーチャービット記号列の総数は 2^N であり、その量はNと共に指数関数的に増加することを注意すべきである。

【0119】

一旦フィーチャーデータ集合が決定されると、どんなサンプルデータ点用にも出力状態確率ベクトルを計算することが出来る。図14を参照すると、このベクトルを計算するためには、全加重係数を創るよう該ローカル及びグローバルエントロピー加重係数を組み合わせることが最初に必要である。本発明の方法では、該ローカル及びグローバルエントロピー加重を含む一般的第3次表現が最適モデル性能用に実験的に調整された係数を用いて規定される。該全加重係数用の一般的表現はかくして次の様に見られる。

【0120】

$$W_{ic}^S = a (W_{i1}^S)^2 W_{c1}^{gs} + b (W_{c1}^{gs})^2 W_{i1}^S + c (W_{i1}^S)^2 + d (W_{c1}^{gs})^2 + e W_{i1}^S W_{c1}^{gs} + f W_{i1}^S + g W_{c1}^{gs} + h$$

かくして、各部分空間S内の各セルiは該与えられた部分空間S用の該ローカル及びグローバル加重の組み合わせである付随する一般的加重係数 W^S を有する(該式は又グローバル加重係数 W_{gs} が出力状態依存性であり、従って該一般的加重係数が出力状態依存性であることを示すことに注意を要す。該グローバル加重係数が全ての出力状態に亘って計算される場合、出力状態cへの依存は除かれる)。

【0121】

aからhまでのパラメーターは最も精密なモデル、フレーム、スーパーフレーム他を得るために実験的に調整される。多くの問題では、該グローバルエントロピー回数も存在するが、該加重係数は該ローカルエントロピー加重係数により支配される。それはここで説明される方法がフィーチャー部分空間内のローカル統計に可成りの重要性を提供する点を強化し、それはここに説明される方法と従来技術のモデル化の取り組みとの間を際立たせる特徴である。該モデル用の信頼限界の確立の中では、該モデル係数は該誤差統計を計算するために変更され得る。

【0122】

一旦 W_{ic}^S 用の適当な値が決定されると、サンプル点d用の各出力状態の確率は次の様に計算出来る。

【0123】

【数15】

$$P_c(d) = \sum_{i=1}^{n_s} W_{ic}^S P_{di}^L,$$

【0124】

ここで該加算は全 n_s 部分空間上に延び、サンプル点dは各部分空間内の対応するセル i_d 内へ射影するよう仮定され、該ローカル確率 $p_{c|i_d}$ は該点がセル i_d 内へ写像する事実がある時、該出力が状態cである確率である。上記の様に、もし一般的エントロピー加重が出力依存でないならば、一般的エントロピー加重の下付き文字cは上記式で無視されてもよい。各出力状態c用確率は次いで確率ベクトル内に組み合わせられ得る。

【0125】

$$P(d) = \{ P_1(d), \dots, P_{K_c}(d) \} / N(i)$$

ここで K_c 出力状態が仮定され、そして

$$N(i) = \sum_{c=1}^{K_c} P_c(i)$$

は正規化係数で、確率の和が1であることを保証するために、 $c = 1$ から K_c までに亘り加算される。

【0126】

10

20

30

40

50

出力状態確率ベクトル $P(i)$ はサンプル点 d の分類までの該データ空間内に含まれた情報を要約している。ニューラルネットワークの様な種々の従来技術のモデル化の取り組みも同様なベクトルとなり、異なる取り組みは該結果を解釈すると取られた。1994年発行の、レビューオブサイエンティフィックインスツルメント (Review of Scientific Instruments)、65巻(6)、1803-1832 pp、ビショップ、シー・エム・(Bishop, C.M.) 著 "ニューラルネットワークとそれらの応用 (Neural networks and Their Applications)" で説明される様に、共通に使用される方法は、予測された出力状態を発生のもっとも大きな確率を有する状態として割り当てる "勝者1人占め (winner take all)" 戦術を使用することである。

フィーチャー部分空間の部分集合を使用する最適モデルの発展

10

高いグローバルエントロピー加重を有する部分空間を同定するための発展型方法は上記で論じられた。これは次元数の災い (curse) が明らかな多くの入力フィーチャーを有する問題で特に有用である。第1の発展段階では、該発展をドライブする適応度関数は部分空間のグローバルエントロピーである。最も良く予測するモデルを決定するために発展の概念を使うことも可能である。第2の発展段階では目標はテストデータ集合で最低誤差となる高いグローバルエントロピーを有するフィーチャー部分空間の最適部分集合を同定することである。この第2の発展段階は最良の予測モデルを作るために協力的仕方で "一緒に良く作用する (work well together)" 部分空間をグループ化する。同時に該モデリング過程で追加的ノイズを導入する部分空間は第2発展段階中に間引かれる (culled)。図15を参照すると、この第2発展段階での該適応度関数は次いで、フィーチャー部分空間の特定の部分集合を使用することから得られるテスト集合内の全体の予測誤差である。

20

【0127】

M が予め決められている第1発展段階の後に M のフィーチャーが高グローバルエントロピーを有するフィーチャー部分空間の最後の遺伝子プール内に存在すれば、フィーチャーの最適組み合わせを見出すために第2発展過程が使用される。 M ビットの "モデルベクトル" が規定されるが、そこでは各ビット位置は与えられたフィーチャーの在り、無しをエンコードする。該モデルベクトルによりエンコードされた該フィーチャーを使用してトレーニングとテストが行われ、該適応度関数はテスト集合上のモデリング過程から生じる適当な性能定量評価である。分類問題用には、該適当な性能定量評価は該テスト集合内に正しく分類されるサンプルのパーセントである。定量的モデリング問題用には、該適当な性能定量評価は該テスト集合内の予測と実際の値の間の正規化された絶対差であり下記で与えられ

30

【0128】

【数16】

$$F = 1 - \frac{\sum_{d=1}^N \frac{|a_d - p_d|}{d_{\max} - d_{\min}}}{N},$$

【0129】

40

ここで a_i はテスト点 d 用の実際出力値、 p_d は該テスト点 d 用の予測値、 d_{\max} はテスト点値の出力範囲の最大値、そして d_{\min} はテスト点値の該範囲の最小出力値である。

【0130】

一旦第2発展過程が終了すると、最適モデルベクトルが該モデリング過程用の最適フィーチャー組み合わせを選択するため使用される。それで、第1発展段階は高情報エントロピーのフィーチャーのプールを同定したが、該プールはテスト集合内の予測誤差を最小にする最良部分集合のフィーチャーを見出すために該第2発展段階で更に発展させられる。この全体の過程は該モデリング問題への最良の実験的解を見出すために種々の発展的条件と制限下で繰り返される。

【0131】

50

かくして本発明の方法は階層的発展の概念を組み入れるが、そこでは最も情報豊富なフィーチャーのみならず、最良予測モデルを開発するために必要なフィーチャー部分空間の最適部分集合も、双方を同定するために、発展的方法が使用される。2つに発展段階を有することは該方法のユニークな利点を提供する。第1段階は手元の問題に見通しを得るために何れの次のモデリング過程からも独立して調べ得るフィーチャー部分空間の情報豊富な部分集合を作る。この見通しは今度は意志決定過程を導くため使用出来る。

【0132】

従来技術のモデリングパラダイムでの共通の苦言はそれらが入力フィーチャー内の何処に情報があるかを容易には明らかにしないことである。この欠点は従来技術の方法の能力を戦略計画と意志決定に参画することを制限する。本発明の方法では、第1発展段階の後の区切り点が、知的戦略計画と意志決定の可能性のみならず、次のモデリング過程が進める価値があるかどうかを決定する機会も考慮する。例えば、もし入力フィーチャーの充分豊富な集合が見出せないならば、本発明の方法は、ローバストなモデルを開発する前に、より情報豊富なフィーチャーを入力として含むデータへ戻るようモデル作成者(modeler)に指し示す。本方法はどの情報がないかを指定はしないが、本方法は充たされる必要のある情報ギャップがあることを指示する。情報ギャップ自体のこの指示は複雑な過程の理解で非常に価値がある。

情報写像の創生(Creation of Information Map)

図11を参照すると、該第1発展段階の後、該問題の基本的理解を得るために該発展したフィーチャーデータ集合内に存在する入力の発生頻度のヒストグラムを作ることにも又非常に有用である。このヒストグラムは該問題用の”情報写像(Information Map)”と規定出来る。幾つかの問題用には、該情報写像の構造は、入力の或る部分集合が入力の他の部分集合より可成り頻繁に起こるならば該問題の次元数を減らすために使用出来る。該部分集合の次元数を減らすことは、セル当たりサンプル点の平均数で部分空間を占めるために必要なデータ量が該次元数の増加につれて指数関数的に増加する様な次元数の災いのもう1つの側面を緩和する追加的利点を有する。図12は遺伝子リストとその付随情報写像の例である。

エグゾースチブ(Exhaustive)な次元的モデリング

図13を参照すると、もしこの様な次元数削減が可能なら、予測モデルは減少した入力データ集合を使用して開発可能である。本方法の好ましい実施例に依れば、Nの最も共通に起こる入力に該情報写像から同定され、次いでNより小さいか等しい全てのM用に該NのフィーチャーのMの部分次元(sub-dimensions)内への全ての起こり得る射影(projection)が該フィーチャー部分空間を規定するため計算される。全てのこの様な射影を計算する帰納的アルゴリズム(recursive algorithm)は次の様である。

【0133】

フィーチャーの全ての組み合わせを計算する帰納的技術(recursive technique)は：各部分次元M用に、Nの数のリスト内で全てのMケ組のもの(M-tuples)(長さMの組み合わせ)を同定する問題を考える。第1要素が最初に選択され次いでN-1の数の残りのリスト内の全ての(M-1)ケ組のもの(長さM-1の組み合わせ)が帰納的仕方で同定される必要がある。一旦全てのこの様な(M-1)ケ組のものが同定され、該第1要素と組み合わせされると、元のリストの第2要素が新しい第1要素として選択され、次いで該第2要素の過ぎた該N-2の残りの要素内の全ての(M-1)ケ組のものが同定される。この過程は該第1要素が該元のリストの終わりからのM+1番目の要素を越えるまで続く。該アルゴリズムはそれがそれ自身を呼ぶので本質的に帰納的であり、それは又該要素の順序付けが重要でないことを仮定している。

【0134】

一旦与えられた部分次元M用の全てのフィーチャーの部分空間のプールが同定されると、このプールは、上記説明の方法を使用してテスト集合内の出力値を予測するために使用されるフィーチャー部分空間の集合として直接使用され得る。この過程は各部分次元M用の複数の量子化条件に亘って繰り返され得る。次いで最適な(部分次元、量子化)-対{op

timum (sub-dimension, quantization) -pairs} がテスト集合上の全予測誤差を最小化することに基づいて選択される。最適な (部分次元、量子化) 対が選択された後、該最適な (部分次元、量子化) 条件に対応するフィーチャー部分空間のプールは該第 2 の発展段階用のスタート点として使用され得る。この第 2 発展段階はテスト集合内に最小全予測誤差を有するこのプールからフィーチャー部分空間の最適部分集合を選択し、かくして最適モデルを規定する。

【0135】

一般的規則として、テスト集合上で十分な全予測精度をなお保存する比較的低い部分次元表現を決定することが有利と分かった。より低い部分次元で、より高いセル母集団統計が量子化の比較的精細なレベルに於いてさえもなお保持され得て、かくして該モデルの精度

10

【0136】

もし元のデータ集合の次元が非常には高くないなら、エグゾースチブな次元モデリングの方法は元のデータ集合に直接適用され得る。これは高情報エントロピーを有するフィーチャーのプールを同定する第 1 発展過程を行う必要性を取り除く。

定量的モデリング

出力変数の人工的量子化を行うことによる定量的モデリング問題の分類問題への変換はローカル及びグローバルエントロピー係数を計算するために有用である。発生する自然な疑問は元のデータ集合内に存在する精度を如何に最終予測モデル内に保存するかである。これは、もし出力ビン解像度が乏しいセル統計を避けるためデータ集合のサイズにより抑制されるならば、特に重要である。伝統的 분류問題用には、出力変数が起こり得る状態の離散的総体 (ensemble) の 1 つを仮定出来るのみなので該精度問題 (precision issue) は存在しない。

20

【0137】

出力変数の人工的量子化を行う 1 つの利点はローカル及びグローバル情報メジャーの計算が、サンプル点の数から共に独立したカテゴリー又はセル上で加算が行われるシャノンの項に基づくことである。これはサンプル母集団統計を情報コンテンツから分離することを容易化する。定量的モデリング用には、出力変数の人工的量子化は該ローカル及びグローバルエントロピーが同じ方法で計算されることを可能にして、かくしてサンプル母集団統計からの情報メジャーの分離を保持する。

30

【0138】

出力変数量子化を使用してローカル及びグローバル情報メジャーが計算された後、生の出力変数内の精度は最終予測モデル内の精度を回復するため使用され得る。

【0139】

最初に出力値の " スペクトラム " が全ての人工的出力変数カテゴリーに亘ってバランスを取られる。これは、各カテゴリー内の最終母集団が共通の目標値にあるように各出力カテゴリー内の各データ項目を或る尺度係数で有効に複製することにより達成される。典型的共通目標値はデータ点の全数を表す数である。

【0140】

データバランス化の 1 方法が上記で説明されたが、特定状態確率 (state-specific probabilities) はその状態に対応する点の数に基づき正規化される。データを明確に複製することなくデータをバランス化する代わりの取り組みを下記で説明する。ニシの情報エントロピー項の計算は、N がデータ集合のサイズを表す場合の $\ln(1/N)$ 係数を含む正規化項を有するが、この正規化は主にエントロピー項を 0 と 1 の間の値に制限するため役立っている。該正規化項は、均一性の程度が該データ集合のサイズに依存する問題に直接向けられていない。

40

【0141】

小さなデータ集合用には、該データ項目の該データ集合内の全データ項目の全体への正規化は微妙な偏倚を招く。例えばデータ内の絶対的変動が比肩されるものでも、より小さいデータ集合内の正規化されたデータ項目間の相対変動は、より大きなデータ集合内の対応す

50

る項目間のそれより大きくなり得る。この偏倚を正すために、データバランス化過程が導入される。該バランス化過程を下記に説明する。

【0142】

2つのデータ集合 D_1 と D_2 を考えるが、ここで該集合はそれぞれ、第1及び第2出力状態に対応する入力を表す。 D_1 は N_1 項目を有し、 D_2 は N_2 項目を有する。 M が N_1 と N_2 の最小公倍数を、 M_1 と M_2 が対応するデータ集合の各々用の掛け算尺度係数 (multiplying scale factors) を表す。もし D_1 を M_1 倍、そして D_2 を M_2 倍だけ複製するなら、最終両データ集合 D'_1 と D'_2 は M 項目を有する。必要な代数計算を行った後、新データ集合の各々用のニシのエントロピー項は次の様に変型される。

【0143】

$$E'_1 = \{ \ln(1/M_1) + f_i \ln f_i \} / \{ \ln(1/M_1) + \ln(1/N_1) \}$$

$$E'_2 = \{ \ln(1/M_2) + f'_i \ln f'_i \} / \{ \ln(1/M_2) + \ln(1/N_2) \}$$

ここで f_i と f'_i はそれぞれ元のデータ集合 D_1 と D_2 上で正規化されたデータ部分を表す。

【0144】

もしセル内の出力データが密にクラスターされていれば、 W_{local} は高い。逆に、もし該出力データが該セル内で全ての人工的出力カテゴリー上にばらまかれていれば、 W_{local} は低い。該グローバルエントロピーは簡単に該部分空間内のセル上での数加重平均 $\langle W_{local}^i \rangle$ として規定出来る。 W_{global} は該部分空間内の情報の正規化総量を測定する。最後に、カテゴリーベースの分類で使用する基本確率定量評価 P^{s_i} は平均(又は代わりに中央値又は他の代表的統計量)セルアナログ出力値で置き換えられ得る。該部分空間上での平均セルアナログ出力値の加重和は次いで出力値を予測する離散的な場合に於ける様に行われることも出来る。それらの出力値で広いばらつき (spread) を有するセルは、個別セルが情報豊富でない部分空間でそうなる様に、下げて加重されることを注意する。

【0145】

セルの平均出力値 μ^s_i の見積もりで、上記で規定したデータ複製尺度係数がバランス化されたデータ集合用にセル内平均値を計算するため使用される。該データバランス化過程はトレーニングデータ集合内の出力値の分布により導入される何等かの偏倚を除去するために行われる。

【0146】

【数17】

$$\mu_i^s = \frac{\sum_{j=1}^n o_j M_j}{\sum_{j=1}^n M_j},$$

【0147】

ここで n はセル内の項目の全数を表し、 o_j は第 j 番の項目の出力値を表しそして M_j は第 j 番のデータ項目に付随するデータ複製係数 (data replication factor) を表すが、該データ複製係数は該第 j 番の項目が属する人工的に量子化された状態に依存する。

【0148】

情報が貧乏なセル及び部分空間からの "クリープ誤差 (creep error)" を減らすために、オプションとして下記の過程が行われる。最初に、情報豊富な部分空間が離散出力状態の議論で前に説明した様に発展させられる。一旦最も情報豊富な部分空間が発展させられると、ローカル及びグローバル両エントロピーしきい値が、該情報豊富な部分空間に付随する平均値か又は中間値か何れかのエントロピー加重和の計算に向かって適用される。該

10

20

30

40

50

ローカルエントロピーしきい値より低いセル用ローカルエントロピー値はゼロ（０）に設定される。同様に、該平均の計算で誤差が徐々に累積されるのを避けるために、該グローバルエントロピーしきい値より低い部分空間用グローバルエントロピー値はゼロ（０）に設定される。

【０１４９】

該ローカル及びグローバルエントロピー関数のしきい値処理（thresholding）で、グローバルエントロピー関数の値の基づき該ローカルエントロピーの追加的しきい値処理を行うことが望ましいことが屢々ある。与えられた部分空間射影用のグローバルエントロピーがその対応するしきい値の下にあれば、その部分空間内の全てのセル用の該ローカルエントロピー関数はそれらの個別値に関係なくオプション的にゼロに設定出来る。前記説明のしきい値処理方法は又離散型出力状態モデリング用にもオプションとして行い得るが、クリープ誤差を最小化するためにより制限的過程が取られるべき定量的モデリング用でより高い価値がある。

10

【０１５０】

最後に、該しきい値処理過程を有しても有さなくても、本発明の方法はサンプルのテスト集合上で最小全出力誤差に帰着する情報豊富な部分空間の最適組み合わせを発展させ得る。又本発明の範囲内の定量的モデリングの方法は階層的発展をも含む。第１発展段階で、最も情報豊富な部分空間が、グローバルエントロピーを適応度関数として使用して、発展させられ、第２発展段階が続くがそこでは最小テスト誤差に帰着する情報豊富な部分空間の最適組み合わせが発展させられる。

20

【０１５１】

従来技術の方法に対する本発明の方法の利点はカテゴリー的及び定量的の両モデリングに共通のパラダイムが使用されることである。実験型のモデリングと過程理解とのための基礎としての分布状階層的発展の概念は、出力変数の唯１つ（連続型か離散型か何れか）の種類用にしか最適化されない従来技術の方法と対照的に、出力変数の両クラス（連続型及び離散型の両方）に適用される。

分布状階層的発展

ここに説明される方法は、“対象（object）”、例えば、フィーチャー、モデル、フレームワーク、そしてスーパーフレームワーク、の階層を創るために、情報理論からの概念を用いて、データの画像的表現、又はデータの多次元的表現の概念を使用する。用語“分布状階層的発展（distributed hierachial evolution）”は、モデル、フレームワーク、スーパーフレームワーク他の様な逐次より複雑で相互作用する発展型“対象”のグループが複雑なデータの漸進的により大きい量をモデル化し理解するため創られる発展型過程として規定される。大きな、複雑なデータ集合用には、前に説明したモデル創生過程が、最適モデルのグループを見出すために種々のトレーニング及びデータ集合上で繰り返される。最適モデルのグループの情報豊富な部分集合は次の様に決定される。

30

【０１５２】

図１６を参照すると、テストデータ集合の入力がモデルの選択された部分集合グループ（ランダムに選択されてよい）の各モデルに差し出され、各部分集合で予測される出力が各テストデータ出力と比較される。該部分集合で予測される出力の計算の過程は個別モデルを創るための過程と同様な仕方で行われ、そこでは個別のモデルで予測される値を入力としてそして実際の出力値を該出力として使用して、新しいトレーニング及びテストのデータ集合が創られる。この過程はモデルの多数の選択された部分集合グループ用に繰り返される。次いで該選択された部分集合グループは、“フレームワーク”と呼ばれるものを規定するためにシステム入力からシステム出力を最も精確に予測するモデルの最適部分集合グループを見出すために発展させられる。図１７Ａと１７Ｂはフレームワーク発展の概念を図解する。

40

【０１５３】

図１８Ａを参照すると、該フレームワーク創生過程は更に、最適フレームワークのグループを見出すためにモデル創生過程と同様な仕方、繰り返される。最適フレームワークの

50

グループの情報豊富な部分集合は次の様に決定される。テストデータ集合の入力がフレームワークの選択された部分集合グループの各フレームワークに印加され、各フレームワーク部分集合で予測される出力が各テストデータ出力と比較される。フレームワーク部分集合で予測される出力を計算する過程は個別モデルを創る過程と同様な仕方で行われるが、そこでは新しいトレーニング及びテストのデータ集合が個別のフレームワークで予測された値を入力として、そして実際の出力値を該出力として使用して創られる。この過程はフレームワークの多数の選択された部分集合グループ用に繰り返される。該選択された部分集合グループは次いで、システム入力からシステム出力を最も精確に予測するフレームワークの最適部分集合グループ（これは「スーパーフレームワーク」と呼ばれる）を見出すために発展させられる。図 18 B はスーパーフレームワーク発展用の考慮を図解する。

10

【0154】

最適モデル決定過程、最適フレームワーク決定過程、或いは最適スーパーフレームワーク決定過程は、予め決められた停止条件が達成されるまで、繰り返されてもよい。該停止条件は、例えば、：1) 予め決められた予測精度の達成、又は2) 予測精度で更に進む改善が達成されない時、の様に規定されてもよい。本発明の方法はかくして実験データ集合上に分布した多数の相互作用する発展型対象の階層が同定される伸長可能な発展型過程である。発展対象の該階層の深さは解析されるべきデータ集合の複雑さにより決定される。簡単なデータ集合用には、全データ集合の非常に小さな部分集合を使用する1つのコンパクトなモデルで該全データ集合に亘りテストと検証 (verification) のデータ集合値を精確に予測するのに充分である。該データ集合の複雑性が増加すると、該全データ集合（検証データ集合を含めて）を精確に説明するためにモデル、フレームワーク、スーパーフレームワークの階層を展開することが必要になるかも知れない。

20

【0155】

分布状階層的発展 (Distributed Hierarchical Evolution) の顕著な計算的利点は、1つの大きな、モノリシックな実験型モデル (monolithic empirical model) の創生よりむしろ実験的モデルを規定するために大きなデータ集合に亘り分布された多数の、コンパクトな発展型対象の創生から生じる。高度に非線形の過程用には、大きなタスクを多くの小さいタスクに分けることが重要な実際の結果を有する顕著な計算的利点を提供する。

【0156】

分布状階層が成長すると、更に最適化が各段階で行われ、全体のデータ集合上での1つの、グローバル最適化上での顕著な性能改善となることは注意されるべきである。該大きなデータ集合内に含まれる益々増える情報は次々とより複雑な発展対象の相互作用の中に閉じ込められ、該相互作用は該実験型モデリング過程内の自由度の顕著な源として作用する。これは新データが現れた時該実験型モデルの更新を簡単化する。該実験型モデルの更新の初期過程は、該新データをテスト集合として使用して現在の実験型モデル内に最も最近の又は「最も高い」発展型対象の新グループを発展させることを含む。より早期のデータを使用して発展させられたより早期の又は「より低い」発展型対象は全く変えられる必要はないが該階層内の最も最近の発展型対象の新グループを創るため使用され得る。より早期の発展型対象のこのリクラスタリング (reclustering) からもし不十分に精確な新実験型モデルが生じるならば、その場合だけ、該新データの部分集合を使用して該階層内の該より早期の発展型対象を再発展 (re-evolve) (該発展の繰り返し) させる必要がある。これが達成された時、最も最近の発展型対象の次ぎに新しいグループが該新データの異なる部分集合を使用して再発展させられる。モデル更新へのこのトップダウン的取り組みは、大抵の従来技術のモデリングの取り組みに共通なより伝統的なボトムアップのモデル更新に勝る顕著な計算的利点を供する。

30

40

監視されないフィーチャークラスタリング

部分集合用グローバルエントロピーメジャーの概念は又入力相関に基づいてフィーチャークラスタを発展させるために適応度関数として使用される。例えフィーチャー部分集合内のセルが出力状態に関し可成りの情報を含まなくても、該セル母集団統計は該部分空間上でなお高度にクラスタされ得る。入力フィーチャー間の相関は、「グローバルエントロ

50

ピー加重係数の代替的規定”の名称の節で前に説明したグローバルエントロピーパラメータの代替の規定と非常に似た情報エントロピー規定を使用して、出力状態から独立にセル母集団統計の均一性を計算することにより同定され得る。この場合、情報エントロピーを計算するために使用されたニシのデータ集合内の基本量はセル母集団であり、該ニシのデータ集合内のエントリーの数該部分空間内のセルの数である。

【0157】

セル占有統計のグローバルエントロピーによりドライブされる発展型技術を使用して、最も高くクラスターされたフィーチャー部分空間は発展させられ、図19A、19B、19Cそして19Dで示される。(19A及び19Bの発展過程は図5A及び5Bの前に説明した過程と同様である。考慮下の特定の遺伝子が過程700で選択される。過程740により示す様に、次の遺伝子シーケンスは過程700で始めに作動させられる。)

これは、クラスターを発見するための、1990年発行、アイイーイーイー論文集(Proceedings of the IEEE)78巻4号1464-1480頁、コーネン、テー。(Kohnen, T.)著”自己組織化写像(The Self-Organizing Map)”で説明される様に、コーネンニューラルネットワーク(Kohnen neural networks)の様な他の監視されない方法の代替である。この様な従来技術の方法に勝る本発明の方法の魅力的側面は監視されない及び監視されるモデリングの間の区別が、該エントロピー計算での出力状態情報の簡単な排除又は包含により非常に自然に起こることである。

【0158】

一旦高度にクラスターされたフィーチャー部分空間のプールが発展させられると、このプール内のフィーチャー部分空間のグループは、帰納用のドライブ条件としての該部分空間を横切る入力重なり用に、例えば、しきい値条件を使用してより大きなクラスターを作るよう帰納的に合併させられ得る。この方法で、より大きなフィーチャークラスターのより小さなグループは、より大きなフィーチャークラスターの直接の同定が計算的に手に負えない非常に高い次元のデータ集合に於いても、効率良く同定され得る。

情報可視化

高いグローバル情報エントロピーのフィーチャーデータ集合を決定する第1の発展段階中に、該発展過程で同定される、最も高いローカル情報エントロピーを有するセルのリストを保持することも又可能である。

【0159】

乏しい、すなわち、人工的に情報豊富なセルのエントリーを避けるためにこのリストの選択では最小セルカウントしきい値が使用されてもよい。高いグローバル情報を有するフィーチャー内に存在するセルを調べることにより第1の発展段階の終わりでこの高いローカルエントロピーリストを創ることは可能である。計算効率の理由で、該第1発展段階の終わりでこの高いローカルエントロピーリストを創ることが好ましい。

【0160】

多次元データ空間内の情報豊富なセルを同定する方法は又”情報可視化(information visualization)”用にも使用出来る。多次元空間での情報可視化はデータ削減の問題として見られる。容易に理解可能な仕方データ集合内の本質的情報を取り込むために、最も情報豊富なセルのみが表示される必要がある。前の段落で、最も情報豊富なセルを選択するシステム的方法が論じられた。一旦これらのセルが全部分空間上で選択されると、カラー科学から得られた方法が視覚的に魅力ある仕方該選択されたセルを表示するため使用されてもよい。例えば、カラー空間の{色相(Hue)、彩度(Saturation)、明度(Lightness)}特徴付けで、該色相座標が該セル出力カテゴリーへ写像され得る。該彩度座標はセルピューリティ(cell purity)のメジャーであるローカルセルエントロピー($E_{L_i}^{L_s}$ か $W_{L_i}^{L_s}$ の何れか)へ写像され得て、該明度座標は該セル内のデータ点の数(すなわち、該母集団)へ写像され得る。他の視覚的写像も行える。該第1発展段階の終わりでカテゴリー当たりのベースで最も情報豊富なセルのアクティブなリストを発生する過程は顕著なデータ減少過程に帰着したことは注意すべきである。このデータ減少は大きなデータ空間内で高い情報のローカル化された定義域(domain)の同定を容易にする。一旦全部分空間上

の走査が該第 1 発展段階の終わりで完了すると、このリストは適当な可視的写像方法を使用して適当な表示装置 { カラーシーアールモニター (color CRT monitor) の様な } 上に表示され得る。かくして多次元データ空間は表示目的で 1 次元リストへ減じられた。本発明の方法のユニークな側面は情報可視化に用いた方法論でデータモデリングを行うため使用された方法論の組み合わせである。両方法用の共通した統合するカーネル (kernel) はセルと部分空間の形式でのデータの画像的表現を用いて情報エントロピーと発展を統合することにある。

ハイブリッドモデリング - 分布状階層的発展のニューラルネットワーク又は他のモデリングパラダイムとの組み合わせ

本方法はデータモデリング用の強力なフレームワークを開示するが、どんなモデリングフレームワークも完全なものはないことを述べることは重要である。全てのモデリング方法は、その取り組み (approach) のためか又は該データに課された構造 (geometries) のためか何れかで、" モデル偏倚 (model bias) " を課す。分布状階層的発展はハイブリッドモデルを創るために他のモデリングパラダイムと組み合わせられ得る。これらの他のパラダイムはニューラルネットワーク又は他の分類又はモデリングフレームワークであり得る。もし他の利用可能なモデリングツールが基本的に異なる哲学を有するなら、それらの 1 つ以上を分布状階層的発展と組み合わせることはモデル偏倚をスムーズ化する効果を有する。加えて、データ偏倚をスムーズ化するために種々のデータ集合を使用して多数の分散されたモデルが各パラダイム内に作られ得る。最後の予測結果は各モデルから来る個別予測の加重された又は加重されない組み合わせとなり得る。かくしてハイブリッドモデリングは、それが種々のモデリング哲学の強さを取り入れるので、極端に強力なフレームワークをモデリングに提供する。

法則の発見 - 分布状階層的発展の遺伝的プログラミングとの組み合わせ

第 1 発展段階の後、生じたフィーチャーデータ集合の情報コンテンツを調べることは教示的 (instructive) である。多くの場合、多数の比較的情報豊富なフィーチャーがあり、それは一緒に用いられると、実験型モデルの次ぎの展開用ベースを形成する。他方、もし、それらの絶対的情報コンテンツ (0 と 1 の間で正規化された) で測定された時、発展させられた情報豊富なフィーチャーがないなら、最も適当な次の過程は、有用でローバストなモデルを発展させるよう努める代わりに該データへ戻ることである。

【 0 1 6 1 】

時々、しかしながら、該第 1 発展段階のもう 1 つの成り行きがあり得る。該データから際立ったフィーチャーが発展することがあるかも知れない。このフィーチャーは極端に情報豊富で、事実、手元の問題用の " 遺伝的コード (genetic code) " を表すかも知れない。この様な場合、より大きなデータ集合が該際立った遺伝子によりコード化された入力を使用して構文解析され得て (can be parsd)、この減少したデータ集合は、下にある法則を説明する数学的表現を発展させるために、遺伝的プログラミングフレームワーク内への入力として使用出来る。遺伝的プログラミングは、例えば、1994 年発行、エムアイテプレス (M.I.T. Pres)、コザ、ジェイ . アール . (Koza, J.R.) 著、" 遺伝的プログラミング - 自然的選択によるコンピュータのプログラミングについて (Genetic Programming - On the Programming of Computers by Natural Selection) " で説明されている。この表現は研究される過程の解析的説明を表し、発展型発見過程の最後の結果である。この過程を用いて、情報理論と発展の組み合わせは、見かけは混乱したシステム内の下にある秩序を閉じ込める数学的表現を発見することに帰着する。情報コンテンツのためにフィーチャーを調べ、次いで実験型モデリングか、数学的発見か、又は該データに戻るか何れかに乗り込む、全体の過程はデータにドライブされるパラダイムに基づく " 発見の科学 (Science of Discovery) " への体系的取り組みを説明する。

【 0 1 6 2 】

混乱したシステムの数学的説明の発展は基本的に内挿的性質 (interpolative nature) か外挿的性質 (extrapolative nature) へと該実験型モデルを変換する。かくして数学的表現は、該実験型モデルの開発で使われるトレーニング集合の範囲の外側でデータ定義域

内に於いてさえ出力値を予測するため使用出来る。又数学的説明はモデル化されつつある過程又はシステム内への基本的見通しと恐らくは下にある原理の発見とを得るための励まし (stimulus) を提供する。

【 0 1 6 3 】

【例】

均質ポリマー連鎖反応 (POLYMER CHAIN REACTION) { ピーシーアール (PCR) } フラグメントの同定

本発明が均質ピーシーアールフラグメントの同定に適用された。本方法は最初にデーヌエイ溶解カーブ (DNA melting curve) の情報豊富な部分を同定し、次いで該入力スペクトラムの情報豊富な部分集合を使用して最適モデルを発展させる。

10

背景

デーヌエイフラグメント同定は伝統的にゲル電気泳動 (gel electrophoresis) により行われて来た。挿入染料 (intercalated dyes) を使用する代替え方法はある時間と感度での利点を提案している。この方法は、加熱時 2 重螺旋デーヌエイが変性する (巻きほぐれる) と該染料蛍光量 (dye fluorescence) が減少することの観察に基づいている。温度に対する蛍光量をプロットする、最終のいわゆる " 溶解曲線 (melt curve) " のデータ解析は該デーヌエイフラグメントのユニークな同定のベースを提供する。しかしながら、該方法は、特定のデーヌエイフラグメントの正確な同定を、他の非特定のフラグメントの存在及び背景基盤 (background matrix) からの蛍光ノイズの存在の両場合で、要求している。

20

スパイク (spiked) される食料サンプルの準備

この研究はピーシーアールを禁ずる知られる食料を評価した。該評価は、該禁止食料の禁止効果を克服するために、該反応へのウシ血清アルブミン (bovine serum albumin) { ビーエスエイ (BSA) } の添加能力をテストした。加えて、溶解曲線解析を使用したピーシーアール製品の均質性検出が臭化エチジウム染色 (ethidium bromide staining) を有する標準的ゲル電気泳動と比較された。

【 0 1 6 4 】

食料は地域の食料雑貨店で購入され、4 で貯蔵された。30 の異なる食料がビーエイエム (BAM) 手順で事前強化 (per-enriched) された。処方された強化法 (enrichment) に従い、サンプルはサルモネラニューポート (Salmonella newport) でスパイクされるか又はスパイクされずに残されたが、表 I I I 参照。該強化は次いでビーエイチアイ (BHI) { デーアイエフシーオー (Difco) } 内で 1 : 10 に薄められ、次いで 37 で 3 時間培養された。

30

【 0 1 6 5 】

【表 2】

表 1

食料	事前強化 ブロス (Broth)	食料 : ブロス 希釈	接種 (Inoculation) レベル
アーモンド	エルビー (LB)	1 : 1 0	0, 10^4 /mL, 10^5 /mL
液状卵	テーエスビー (TSB)	1 : 1 0	0, 10^4 /mL, 10^5 /mL
赤小麦ぬか (Red Wheat Bran)	エルビー	1 : 1 0	0, 10^4 /mL, 10^5 /mL
ピーナッツバター	エルビー	1 : 1 0	0, 10^4 /mL, 10^5 /mL
ウォールナット	エルビー	1 : 1 0	0, 10^4 /mL, 10^5 /mL
挽きコーヒー	エルビー	1 : 1 0	0, 10^7 /ml
インスタントコーヒー	エルビー	1 : 1 0	0, 10^7 /ml
インスタント茶	エルビー	1 : 1 0	0, 10^7 /ml
タイム (Thyme)	テーエスビー	1 : 1 0	10^7 /ml
チョコレートアイスク リーム	脱脂粉乳 (Non fat dry milk)	1 : 1 0	10^7 /ml
バジル	テーエスビー	1 : 1 0	10^7 /ml
ホットチョコレートミ ックス	脱脂粉乳	1 : 1 0	10^7 /ml
オレガノ (Oregano)	テーエスビー	1 : 1 0 0	10^7 /ml
ペーストリナットミッ クス (Pastry Nut Mix)	エルビー	1 : 1 0	10^7 /ml
全スパイス	テーエスビー	1 : 1 0 0	10^7 /ml
ローズマリー	テーエスビー	1 : 1 0	10^7 /ml
シナモン	テーエスビー	1 : 1 0 0	10^7 /ml
小麦ぬか	エルビー	1 : 1 0	10^7 /ml
カーネーション、ホッ トココアミックス	脱脂粉乳	1 : 1 0	0, 10^7 /ml

【 0 1 6 6 】

【表 3】

食料	事前強化 ブロス (Broth)	食料： ブロス 希釈	接種 (Inoculation) レベル
ネスルのココア	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
オレオのクラム (Oreo Crumb)	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
スイスモカコーヒー	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ネスルチョコレトリカー	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ミルクチョコレート	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ハーシーのココア	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ダークココア (Dark Cocoa)	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ウイナチョコレート コーヒー (Viennese Chocolate Cafe)	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ウォールナットホイップ (Walnut Whip)	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml
ネスルのミルクチョコ レートクラム (Nestle's milk chocolate crumbs)	脱脂粉乳	1 : 1 0	0, 10 ⁷ /ml

【 0 1 6 7 】

ポリビニルポリピロリドン (Polyvinylpolypyrrolidone) { ビーバイビーピー (PVPP) } 処理

グロバックサンプル (growback) の 5 0 0 マイクロリットル (500 ul) のアリコート (aliquot) がビーバイビーピー { クアリコン社 (Qualicon, Inc.) } の 5 0 m g のタブレットを含むチューブに追加された。該チューブはボルテックス (vortexed) されそして該ビーバイビーピーは 1 5 分間澄むようにされた。最終浮遊物は次いで溶解過程で使用される。

サルモネラサンプルの準備

2 m l のスクリーカップチューブ (screw cup tube) で、強化すなわちビーバイビーピー処理サンプルの 5 マイクロリットルがデーエヌエイ挿入染料エスワイビーアールグリーン (DNA intercalating dye SYBR^R Green) { モレキュラープローブ (Molecular Probes) } の 1 : 1 0、0 0 0 希釈を含む溶解試薬 { 5 m l ビーエイエックス溶解バッファー (5ml BAX^R lysis buffer) と 6 2 . 5 u l (マイクロリットル) ビーエイエックスプロテアーゼ (62.5 ul BAX^R Protease) } の 2 0 0 u l (マイクロリットル) に加えられた。該チューブは 3 7 ° で 2 0 分間次いで 9 5 ° で 1 0 分間培養された。9 5 ° の培養の後、4 m g / m l のビーエスエイ (BSA) 溶液の 5 0 u l (マイクロリットル) が該溶菌液 (lysate) に追加された。これはビーバイビーピー処理済みと未処理のサンプルに行われた。対照として、幾つかのサンプル未処理で残された。この未精製バクテリア溶菌液の 5 0 マイクロリットルが、パーキンエルマー 7 7 0 0 シークエンスデテクター計器 (Perkin Elmer 7700 Sequence Detector instrument) で使用されるピーシーアールチューブ内に含

まれた1つのビーエイエックスサルモネラサンプルタブレット (BAX^R Salmonella sample tablet) を水和するため使用された。該チューブはキャップを付けられ、パーキンエルマー 9600 サーマルサイクラー (Perkin Elmer 9600 thermal cycler) 内で次のプロトコルに依り熱サイクルにかけられた。

【0168】

94 2 . 0 分 1 サイクル
 94 1 5 秒 3 5 サイクル
 72 3 . 0 分
 72 7 分 1 サイクル
 4 " 長期間 (forever) "

10

増幅後分析 (Post Amplification Analysis)

増幅後、下記条件で運転することによりパーキンエルマー 7700 デーエヌエイシーケンスデテクター (Perkin Elmer 7700 DNA Sequence Detector) 上で該溶解曲線が作られた。

【0169】

プレートの種類： シングルリポーター (Single Reporter)
 器械： 7700 シーケンスデテクションシステム (7700 Sequence Detection System)
 運転： 実時間
 染料層： エフエイエム (FAM)
 サンプルの種類： 未知である
 サンプル容積： 50 u l (マイクロリットル)
 運転条件：

20

70 2 分 1 サイクル データ収集せず
 68 1 0 秒 9 8 サイクル データ収集する
 自動インクレメント + 0 . 3 / サイクル
 25 " 長期間 "

該多成分データは該器械から移出され該分析に使用された。特定のデーエヌエイフラグメントの製作は該アンプリファイ (amplified) されたサンプルにビーエイエックスローディングダイ (BAX^R Loading Dye) の 1.5 マイクロリットルを添加することにより検証された。次いで 1.5 マイクロリットルのアリコートが臭化エチジウムを含む 2 % アガロースゲル (agarose gel) のウェル (well) 内に装填された。該ゲルは 30 分間 180 ボルトで運転された。特定の生成物は次いでユーブイトランシイルミネーション (UV transillumination) を使用して可視化された。

30

データ分析

生の蛍光量 (raw fluorescence) データが処理用にマイクロソフトエクセル (Microsoft Excel) に移入された。この段階からデータを可視化し該データから予測をするため分岐的取り組みが使用された。

データ事前処理 (Data Preprocessing)

蛍光ノイズを減らすために該データを事前処理することは成功するモデリングの尤度 (likelihood) を増すことが実験的に決定された。該データ事前処理は次の過程から成り、すなわち、

40

a . 蛍光データ (fluorescence data) の正規化、
 b . 0 . 1 の解像度でキュービックスプライン関数 (cubic spline function) を用いた該正規化蛍光の内挿補間、
 c . 内挿補間された蛍光スペクトラムの対数を取る、
 d . 25 点サビツスキーゴレイ平滑化関数 (25 point Savitsky Golay smoothing function) を用いた該蛍光の対数の平滑化、
 である。

【0170】

50

最終温度スペクトラムはここで説明されるモデリング方法への入力の実集合として使用される。該温度スペクトラムを使用した2つの異なるモデリング例を説明する。

過程 a . データの正規化と可視化

該蛍光データは、最初にスペクトラム内の最低測定蛍光レベルを決定し、この値を、直流オフセットを除くために、該スペクトラム内の各点から引くことにより正規化される。上記の過程 a . の正規化されたデータは次いでサビツスキーゴレイの平滑化アルゴリズム (Savitzky-Golay smoothing algorithm) で平滑化される。温度に対する平滑化蛍光の負の導関数 $\{-d \log(F) / dT\}$ が取られ、 $-d \log(F) / dT$ (y 軸) 対温度 (x 軸) としてプロットされる。

過程 b . 該データからの予測

該正規化されたデータからスタートして、キュービックスプライン内挿関数 (cubic spline interpolating function) を使用して 0 . 1 C 分解能で該データは内挿補間される。次いで該内挿されたデータの対数が取られ、次いで 2 . 5 度 (すなわち 0 . 1 で 25 の点) 上でサビツスキーゴレイの平滑化アルゴリズムを用いて平滑化される。温度に対する該ログの蛍光の負の導関数 $\{-d(\log F) / dT\}$ 、サルモネラ用データ範囲: 82 . 0 - 93 . 0 (12 データ点) を用いて 1 . 0 C 間隔でパース (parsed) された。

【0171】

方法比較用に、ここに説明された方法は2つの他の良く知られたモデリング方法: ニューラルネットワーク及びロジスティック回帰 (logistic regression)、と比較され、結果は下表で報告される。

【0172】

見出された最も有効な DNA フラグメント同定法は2つのモデリングスキームをシーケンシャルな仕方で背中合わせで使うことを含んでいる。同定の第1レベルはスメア (smear) を非スメア (non-smear) から分離することである。これに、非スメアサンプル用に關心のある特定のデーエヌエイフラグメントを同定することが続く。実際は、この階層的な方法は、起こり得る出力カテゴリーを表す正、負そしてスメアを有する1つの3状態モデルを使用するより精確であった。

1 . 特定ピーシーアールフラグメントに対する非特定ピーシーアールフラグメントのモデリング

該ピーシーアールアンプリフィケーション過程 (PCR amplification process) は、關心のあるデーエヌエイの特定の種類に対応するフラグメントのみならず非特定ピーシーアールフラグメントも作る。第1例は本方法の該非特定と特定のピーシーアールフラグメント間を区別する能力を展示する。149のロックされたプロセス (すなわち、対照) 特定のトレーニングスペクトルと、問題食料 (ピーシーアール用で問題があると知られる実際の食料) の309のテストスペクトルと、一緒に30の非特定の又は " スメア " の蛍光スペクトルのグループが創られた。0 . 1 の温度分解能を有して、111点を含む各サンプル用の温度スペクトル (11 . 1 の範囲上の) が創られた。該ロックされたプロセスと問題食料サンプルの両者が陽性と陰性の標本を含んだ。この例で、該陽性のサンプルは特定のバクテリア (例えば、サルモネラ) でスパイクされ (すなわち汚染され) そして陰性のサンプルはスパイクされぬ (汚染されぬ) ようにされた。該スメアサンプルはロックされたプロセストレーニング集合 (12 スメアサンプル) と問題食料テスト集合 (18 スメアサンプル) の両者にランダムに導入された。該陽性及び陰性の両サンプル状態は合併され2進のゼロ " 0 " 文字でラベル付けされ、該スメアサンプル状態は2進の1 " 1 " でラベル付けされた。

【0173】

a . 入力の最も情報豊富な集合を発展させること

モデリング過程の第1歩は111次元の入力フィーチャー空間をより少ない、より情報豊富な部分集合に減じることである。前に説明した発展型フレームワークが該最も情報豊富なフィーチャーを発展させるために使用された。100の遺伝子の初期遺伝子プールがラ

10

20

30

40

50

ンダムに発生され、そこでは各遺伝子は2進の111ビットの長さの記号列を有し、各ビットの状態は該対応入力フィーチャーが該遺伝子内で賦活されたかどうかを表している。該発展過程はセル当たり1サンプルとなるべき平均セル占有数 (mean cell occupation number) により抑えられ、そして該発展は5世代より多く進んだ。各遺伝子の発展をドライブするために、グローバルエントロピー、又は適応度関数としてローカルエントロピーの数加重和 (number-weighted-sum of local entropies) が使用された。該発展は固定サイズ化された部分範囲 (すなわち、適応型ビンニングよりむしろ、固定されたビン) を使用して進みそして該データは、上記説明の様に、0及び1の出力状態の数をバランスさせるようバランスさせられた。

【0174】

10

発展型過程を通して該100の最も情報豊富な遺伝子のグローバルリストが保持された。全ての111の入力フィーチャーのビット頻度のヒストグラムが、発展した該情報豊富な遺伝子プール内で最も屢々発生するビットを同定するために、該発展の各世代の終わりで分析された。このヒストグラムはどの温度点が該出力状態に最も密接に付随したかについての情報を提供した。

【0175】

該111の点の温度範囲が0から110までインデックス (indexed) され、下記31温度点が該発展型過程から選択された: 12, 14, 16, 18, 20, 22, 24, 26, 28, 30, 32, 34, 36, 38, 40, 42, 44, 46, 50, 52, 54, 56, 58, 60, 62, 64, 80, 82, 84, 86, 88。

20

【0176】

情報豊富な領域が該ヒストグラム内で観察されそしてこれらの領域に懸かる偶数番号インデックス点 (上記リスト) が選択されたことは注意されるべきである。大抵の該選択された点が12から60の範囲に懸かることは注意されるべきである。これは該スメアサンプル用溶解曲線スペクトラムが該ベースライン上に立ち上がりそして該インデックス間隔 [12, 60] に対応する温度範囲内の陽性及び陰性両サンプルから別れ始めるからである。例えスメアがそれらの正に規定により可変溶解曲線構造を有するとは云え、主な構造的フィーチャーは該陽性のサンプル内よりも低い温度で一般に現れる。該陰性のサンプルは本質的に構造から自由である。かくして、本方法はより低い温度領域がスメアと非スメアの間の最良の区別が起こる場所であることを確認する。

30

【0177】

b. パース (parsed) されたデータの全低次元射影のエグゾーストな探索
第1発展型過程で発見された該情報豊富な点を使って該トレーニングデータ集合がパースされた後、該減少したデータ集合は広いビンニング範囲に亘り低次元でエグゾースチブに探索された。固定ピンとデータ集合パラランシングが該エグゾースチブな過程を通して使用された。このモデリング問題で、次元当たり26の固定ピンを使用して全2次元射影内への該31次元入力空間の465の射影を発生することが該最良エグゾースチブモデルに帰着することが分かった。 $W_1^2 = 10$ 、 $W_1 = 5$ 、定数項 = 1のエントロピー加重係数が使用された。しかしながら、全465の射影を使用する該エグゾースチブモデルは、該射影の多くが情報より多くのノイズを導入するので、最適モデルであることを保証されない。それで、各ビットが該モデル用遺伝子プール内の与えられた2次元射影の包含 (inclusion) (2進で1) と排除 (exclusion) (2進で0) を表す465ビットの長さの2進記号列を使って第2の発展段階が行われた。

40

【0178】

c. 最良2次元モデルを発展させること
100のランダム2進記号列が最初に発生されそしてそれらの適応度関数がテストデータ集合内誤差を該発展型過程をドライブする適応度関数として使用して計算された。該モデルは20世代より多く発展させられそして最も情報豊富な遺伝子のグローバルリストが保持された。最後に、この遺伝子プール内の最も情報豊富な遺伝子 (最小テスト誤差に帰着する遺伝子に対応する) がスメア検出用遺伝子コードとして選択された。この遺伝子は

50

該包含 2 次元射影の 1 6 3 を有し残りの射影は排除された。これらの 1 6 3 の射影を使用した最小テスト誤差は該 3 2 7 テストケースから 3 つのエラー (3 errors out of the 32 7 test cases) (3 0 9 問題食料サンプルと 1 8 スメアサンプル) であって 9 9 % より高いモデル精度に帰着する！

2 . 陰性のサンプルに対する特定のサルモネラピーシーアールフラグメント (陽性の) のモデリング

ピーシーアールモデリングの第 2 例として、本方法は食料サンプル内サルモネラに対応する特定のデーエヌエイフラグメントを同定するタスクを与えられた。もう 1 度、該ロックされた過程スペクトルが該トレーニングデータ集合として使用されそして該問題食料スペクトルが該テストデータ集合として使用された。上記説明のものと同様な過程が最良予測モデルを発展させるために使用された。

10

【 0 1 7 9 】

a . 入力のも最情報豊富な集合を発展させること

前の例で説明されたそれと同様な手順に従い、本方法は、下記の温度点：

1 0 , 1 3 , 1 6 , 6 1 , 6 4 , 6 7 , 7 6 , 7 9 , 8 2 , 8 5 , 8 8 , 9 1

に対応する 1 2 入力フィーチャーの集合を発展させた。

【 0 1 8 0 】

この例では、スペクトルの情報豊富な部分は該温度範囲のより高い端 (点 6 1 から 9 1 の間) 内にあることを注意する。これは余り驚くべきことではないが、それはポジティブな (positive) 溶解曲線内の主な構造が温度インデックス (temperature index) 8 0 の周

20

辺で起こるからである。

【 0 1 8 1 】

b . パースされたデータの全低次元射影のエグゾースチブな探索

第 1 発展過程で発見された該情報豊富な点を使用して該トレーニングデータ集合がパースされた後、減少したデータ集合は広いビニング範囲上で低次元でエグゾースチブに探索された。固定ピンとデータ集合バランスングが該エグゾースチブな過程を通して使用された。このモデリング問題で、次元当たり 1 9 の固定ピンを使用した全 3 次元射影内への該 1 2 次元入力空間の 2 2 0 の射影を発生することが最良エグゾースチブモデルに帰着することが分かった。前のサンプルで同じエントロピー加重係数が使用された。この例で、全ての 2 2 0 の射影を使用することが最良モデルに帰着することが分かった。該 2 2 0 の射影の部分集合を発展させることは該テストデータ集合に関する予測精度を改良しなかった。全 2 2 0 の射影を用いて、該 3 0 9 の問題食料テストサンプル (スメアなしで) からの 3 0 1 が 9 7 . 4 % の精度で適当と同定された。

30

結果

これらの実験中作られた該 3 0 9 のデータサンプルの中で、2 0 4 はサルモネラでスパイクされそして 1 0 5 のサンプルが " ブランク (blank) " 反応であった。該 2 0 4 のスパイクされたサンプルの中で、1 4 3 のサンプルはアガロースゲルで陽性でありそして 6 1 は該ゲルで陰性であった。該陰性のサンプルはピーシーアールの禁止か又は不適当なゲルか又はピーシーアール感度の結果と考えられ得る。該 1 0 5 の " ブランク " の反応の中で、9 5 は該ゲルに関し陰性で、そして 1 0 は該ゲルに関し陽性であった。該陽性のサンプルは自然の食料汚染 (例えば、液状卵サンプル) 又は技術的誤りの結果と考えられ得る。

40

【 0 1 8 2 】

下表は該 3 つのモデリング方法の結果を抄録する。該モデリング方法の各々の出力は 1 かゼロの間の数である。 " 1 " はスパイクされた予測を表す一方 " 0 " はスパイクされていない予測を表す。該数がゼロ又は 1 に近い程、該予測により高い信頼を置くことが出来る。0 . 5 のしきい値より高いどんな予測も陽性と考えられた。下記方法の各々用数は期待予測と合致したサンプル数を示す

【 0 1 8 3 】

【表 4】

表 I I

	説明	期待される予測 ²	サンプル数 ³	本方法	ニューラルネット	ロジスティック回帰
スパイクされた／陽性ゲル	確認された陽性	1	143	139	138	134
スパイクされぬ／陰性ゲル	確認された陰性	0	95	93	92	64
スパイクされぬ／陽性ゲル	被汚染サンプル	1	10	8	8	10
スパイクされた／陰性ゲル ¹	検出感度	0／1	61	56／5	55／6	47／14
合計			309	301	299	269
			%合致度	97.4 1%	96.7 6%	87.0 6%

【0184】

¹これらのサンプルはスパイクされたが、ゲル上では陰性であった。均質な検出はゲル検出より敏感なので、均質な検出で陽性のサンプルを検出するがゲルベースの方法では見出さないことが起こり得る。パーセント合致度計算時、このカテゴリーで全てのサンプルは正しいと仮定されている。

² "期待される予測" 列はスパイクステータスとゲル結果とに基づき1又は0を表示する。この数は該モデルが該トレーニングサンプルに基づき予測すると期待されたものである。

³ "サンプル数" 列は特定のスパイク／ゲルカテゴリーに分類されるサンプル数を表示する。

【0185】

本方法の階層化モデリングに加えて、ハイブリッドモデリングフレームワークが使われてもよい。

【0186】

ニューラルネットモデルは陽性／陰性の同定のみならずスミア／非スミアの同定用にも開発された。事実、より多くのデータが入手可能になると、多数のトレーニング／テストデータ集合が発生され得て多数ニューラルネット及びインフォエボリューションモデル (InfoEvolve™ model) に帰着した。未知のサンプルは全てのモデルでテストされ得て個別モデル予測の統計に基づきカテゴリー化され得る。付録Gで論じる様に、この取り組みは、多数のデータ集合とモデリングパラダイムと上での多様化によりモデル偏倚のみならずデータ偏倚も減じる利点を有する。加えて、2つの別々のモデリング段階を続けて使用する

階層的取り組みはモデル精度を更に改善する。

ハイブリッドモデリング

本方法はデータモデリング用の強力なフレームワークを開示するが、どんなモデリングフレームワークも完全ではないことを注意することは大切である。全てのモデリング方法はその取り組みのためか又はデータに課されるジオメトリー (geometries) のためか何れかで、"モデル偏倚" を課す。本方法は追加的ジオメトリーの最小の使用を行いそして上記説明の様に幾つかの利点を有するが、しかしながら、本方法は基本的に外挿法的であるより寧ろ内挿法的である。比較的データの貧弱なシステムでは、この内挿法的特性は一般化の容易さを減じる。

【0187】

本方法の強さを利用しそしてその弱さを最小化するために、それはハイブリッドモデルを創るために他のモデリングパラダイムと組み合わせられることが可能である。これらの他のパラダイムはニューラルネットワーク又は他の分類又はモデリングフレームワークであり得る。もし他のモデリングツール (含む複数ツール) が基本的に異なる哲学を有するなら、1つ以上の他のモデリングツール (含む複数ツール) を本方法と組み合わせることがモデル偏倚を平滑化する (smooth out) 効果を有する。加えて、データ偏倚を平滑化するために異なるデータ集合を使用して各パラダイム内に多数のモデルが作られ得る。最後の予測結果は各モデルから来る個別予測の加重又は非加重の組み合わせとすることが出来る。ハイブリッドモデリングは多様なモデリング哲学の強さを利用するために極端に強力なフレームワークをモデリングに提供する。重要な意味で、この取り組みは実験型モデリングの究極の目標を表す。

【0188】

例えば、もし食料媒介病原菌用テスト (testing for foodborne pathogens) での上記説明例に於ける様に、偽陰性のパーセント (percento of false negative) を最小化したい望みがあるなら、該モデルのどれか1つがスパイクされたサンプルを予測したならば陽性の結果が報告されるであろう。もしこの規則がこの例のデータに適用されたなら、ゲル結果に基づく偽陽性 (false positive) の率は0.7%より少なかったであろう。何れか1つのモデルについての偽陰性率はそれぞれ: 本方法 = 3.9%、ニューラルネットワーク = 4.5%そしてロジスチック回帰 = 5.8%であった。

結論

この例は重要な実験型モデリング問題でのインフォエボルブテーム (InfoEvolve™) のパワーを図解する。インフォエボルブテームは最初にデーエヌエイ溶解曲線の情報豊富な部分を同定し次いで該入力スペクトラムの情報豊富な部分集合を使用して最適モデルを発展させる。この例で追跡された一般的パラダイムは種々の産業及びビジネス応用品でテストされ大きな成功をもたらし、この新しい発見的フレームワークに強力な支持を提供している。

製造過程の例

ケルバーオール (Kelvar^R) 製造過程での重要な変数は該ケルバーオールパルプ (Kelvar^R pulp) 内に保持された残留湿気 (residual moisture) である。該保持された湿気は該パルプの次の処理可能性と最終製品特性の両者に顕著な影響を有する。かくして最適制御戦略を規定するために該パルプ内の湿気保持に影響するキー要素、又はシステム入力を最初に同定することが重要である。製造システム過程は、乾燥処理用の全体の時間枠のために該入力変数と最終パルプ湿気間の多数の時間遅れの存在により複雑化される。パルプ乾燥処理のスプレッドシートモデルが創られ得るが、そこでは該入力はいくつの前の時の幾つかの温度と機械的変数を表し、該出力変数は現在時刻のパルプ湿気である。最も情報豊富なフィーチャー組み合わせ (又は遺伝子) は、その変数の、より早期の時点でパルプ湿気に影響するのに最も情報豊富であるのはどの変数であるかを発見するためにここに説明された該インフォエボルブテーム (InfoEvolve™) を使用して発展させられ得る。

フロード (fraud) 検出例

既知のフロード的 (fraudulent) な場合のトレーニング集合を作るのが難しいからだけで

10

20

30

40

50

なく、フロードが多くの形式を取るかも知れないので、フロード検出は特に挑戦的応用である。フロードの検出は予測モデリングによりフロードを防止出来るビジネス用に可成りのコスト節約へ導き得る。フロードが起こる或るしきい値確率で決定出来る様なシステム入力の同定が望ましい。例えば、何が " ノーマル (normal) " な記録かを最初に決定することにより、或るしきい値より多く該ノーム (norm) から変化する記録が、より精密な精査用にフラグ建て (flagged) されてもよい。これは、クラスタリングアルゴリズムを適用し、次いでどのクラスターにも分類されない記録を調べることに依るか、又は各分野の値の期待範囲を説明する規則を作ることに依るか、又は分野の異常な付随にフラグ建てすることにより行われてもよい。クレジット会社は期待しない使用量パターン (usage patterns) にフラグを建てるこのフィーチャーをそれらの課金正式化過程内にルーチ的に組み込む。もしカード所有者 (cardholder) が普通は彼 / 彼女のカードを航空券、レンタルカー、そしてレストラン用に使用するが、或る日それをステレオ機器か又は宝石を買うため使用するなら、その処理は、該カード所有者が彼のアイデンティティを検証する該カード発行会社の代表者と話を出来るまで、遅延してもよい。(参考文献: 1997年発行、マイケル、ジェイ・エイ・ベリー、及びゴードン、リンホフ (Michael J. A. Berry, and Gordon Linhoff) 著、" マーケティング、販売及び顧客サポート用データマイニング技術 (Data Mining Techniques for Marketing, Sales, and customer Support)、76 ページ)。フロード検出でどの変数が最も情報豊富かを発見するために最も情報豊富なフィーチャー組み合わせ (又は遺伝子) がここで説明した本発明を使用して発展させられ得る。これらの変数は或る時間間隔に亘る購入の種類と量、クレジットバランス、最近の住所変更他を含んでもよい。一旦入力の情報豊富な集合が同定されると、これらの入力を使用する実験型モデルは本発明を使用して発展させられ得る。これらのモデルは、フロード検出用の適合学習型フレームワークを創るために、新データが入ると規則的ベースで更新され得る。

マーケティング例

銀行は予防的アクションを行う時間を持つためにその要求払い預金勘定 (demand deposit accounts) { 例えば、銀行当座預金 (checking accounts) } の顧客のアトリッション (attrition) の十分な警報を望む。それが余りに遅くなる前にトラブル範囲に見つけるために、起こり得る顧客のアトリッションをタイムリーな仕方で予測するキー要素又はシステム入力を決定することが重要である。かくして、勘定動向 (account activity) の毎月の抄録はこの様なタイムリーな出力を提供しないが、処理レベルでの詳細データは提供するかも知れない。システム入力は、顧客が該銀行に置いて行く理由を含んでおり、この様な理由がもっともかどうかを決定するためにデータ源を同定し、次いで該データ源を処理経過データと組み合わせる。例えば、顧客の死亡が処理停止の出力を提供したり、或いは顧客は最早2週間毎に支払われないか又は最早直接預金を有せずかくして規則的な2週間ベースの直接預金は最早ない。しかしながら、内部決定で発生されたデータは処理データ内に反映されない。例は、該銀行がかって無料であったデビットカード処理用に今は課金しているから又は該顧客がローンのために拒絶されたから、顧客が去って行くことを含んでいる。{ 1997年発行、マイケル、ジェイ・エイ・ベリー、及びゴードン、リンホフ (Michael J. A. Berry, and Gordon Linhoff) 著、" マーケティング、販売及び顧客サポート用データマイニング技術 (Data Mining Techniques for Marketing, Sales, and Customer Support)、85 ページ参照 }。予測的アトリッションを決定する中でどの変数が最も情報豊富であるかを発見するために、ここで説明した本発明を使用して最も情報豊富なフィーチャー組み合わせ (又は遺伝子) が発展させられ得る。顧客属性のみならず銀行戦略に付随する内部管理も含めた両者が処理データパターンと組み合わせられるデータベースを創ることは銀行戦略、顧客属性そして発見されるべき処理パターンの間の起こり得る情報豊富なリンケージを可能にする。これは今度は処理挙動を予測する顧客挙動予報モデル (customer behaviour forecasting model) の発展へ導くことが出来る。

金融予測例 (Financial Forecasting Example)

金融予報 { 例えば、株、オプション、ポートフォリオ (portfolio) そして物価指数 (ind

10

20

30

40

50

ex pricing) } での重要な考慮は株式市場の様な動的で移り気な活動場所では誤差の広い
マージンを黙認する出力変数を決めることである。例えば、実際の物価レベルよりむしろ
ダウジョーンズ平均株価指数 (Dow Jones Index) での変化を予測することは誤差のより広
い許容限度 (wider tolerance for error) を有する。一旦有用な出力変数が同定される
と、次の過程は最適予測戦略を規定するために該選択された出力変数に影響するキー要素
、又はシステム入力を同定することである。例えば、ダウジョーンズ平均株価指数の変化は
ダウジョーンズ平均株価指数での前の変化のみならず他に於ける国の及びグローバルの指数
にも依存するかも知れない。加えて、グローバルな利率、外国為替レート及び他のマクロ
経済的メジャー (macroeconomic measures) が重要な役割を演ずる。加えて、最も金融的な
予報問題は入力変数 (例えば、前の価格変化) と終わりのタイムフレームでの最後の価格
変化との間の多数の時間遅れの存在により複雑化する。かくして、該入力はその前の多数の時
刻での市場変数 { 例えば、価格変化、市場の移り気 (volatility of the market) 、移り
気モデルの変化 (change in volatility model) 、 . . . } を表しそして該出力変数は現
在の時刻での該価格変化である。 (参考文献 : 1996 年発行、エドワードゲートレイ (10
Edward Gateley) 著、 " 金融予測用ニューラルネットワーク (Neural Networks for Fina
ncial Forecasting) 、 20 ページ) 。 より早期の時期が指すどの変数が金融予測用市場変
数への影響で最も情報豊富であるかを発見するためにここで説明する本発明を使用して最
も情報豊富なフィーチャー組み合わせ (又は遺伝子) が発展させられ得る。一旦これら (変
数、時点) の組み合わせが発見されると、それらは最適金融予測モデルを発展させるた
めに使用出来る。 20

【 0 1 8 9 】

下記はモデル発生にここで使用される説明した方法に関する擬コードリスティング (Pseu
de Code listing) である :

```

LoadParameters();          //データ集合と、ピニングの種類の様な種
                             々のパラメーターとをロードし、データ選出、
                             エントロピー加重係数、データ部分集合の数
                             他. . . をバランスさせる

Loop through subset_number {
  CreateDashSubset(filename) //部分集合データをランダムに
  Loop through number of local models {
    EvolveFeatures();        //情報豊富な遺伝子を発展させる
    CreateTrainTestSubset(); //データ部分集合をトレイン/テスト部分
                             集合に分ける
    EvolveModel();           //モデルを発展させる
  }
}

CreateDataSubset
  DetermineRangesofInputs;
  if(BalanceStatsPerCatFlag is TRUE)
    BalanceRandomize;
  else
    NaturalRandomize;
DetermineRangeofInputs
  Loop through data records {
    Loop through input features {
      if(input feature value=max
        or input feature value=min {
        LoadMinMaxArray (feature index, feature value) ;
        UpdateMinMax (feature value) ;
      }
    }
  } //入力フィーチャーループ終了
} //データループ終了

BalanceRandomize
/ * * * * *
/ データ集合を現在の部分集合と残りの部分集合とに分ける ;
/ 出力カテゴリー当たりの項目の数をユーザーが指定する。
/ * * * * *

```

10

20

30

```

Loop through output stats {
    InitializeCountingState(output) to 0;
    InitializeCountingRemainingState(output) to 0;
}
Loop through data records {
    Set IncludeTrainFlag to FALSE;
    Loop through input features {
        if(input features =min) {
            if(input FeatureMinFlag=CLEAR){
                IncludeTrainFlag=TRUE;
                FeatureMaxFlag =SET;
            }
        }
        elseif (input feature=max) {
            if (input FeatureMaxFlag=CLEAR) {
                IncludeTrainFlag=TRUE;
                FeatureMaxFlag =SET;
            }
        }
    }
    //フイーチャーループ終了
    output=ReadOutputState; //記録用に出力状態を読み出す
    guess=GuessRandomvalue;
Threshold(output)=NUMITEMSPERCAT/TotalCountinState (output)
//TotalCoutinState (output) は出力カテ
//ゴリー内の#データ項目を意味する

/ * * * * *
もしデータ記録がフイーチャー最小又は最大値の最初の場合なら、現在のデータ部分集合
と残りのデータ部分集合の両者へ記録をコピーする。
/ * * * * *

    if (IncludeTrainFlag=TRUE) { //現在の部分集合と残りのデータ部
                                分集合の両者へ記録をコピー
        CopyRecordtoCurrentDataSubset;
        IncrementCountinState (output) ;
        CopyRecordtoRemainingDataSubset;
        IncrementCountinRemainingState (output) ;
    }

/ * * * * *
或いは他にもし該出力カテゴリーの項目の数が過剰にNOTであるなら、該データ項目を該R
EMAININGデータ部分集合内に置き換える。
/ * * * * *

```

10

20

30

40

```

elseif(Threshold (output) > MINIMUM_THRESHOLD) {
    CopyRecordtoRemainingData;
    IncrementCountinRemainingState (output) ;
    if (CountinState(output) < NUMITEMSPERCAT) {
        CopyRecordtoDataSubset;
        IncrementCountinState (output) ;
    }
}

//MINIMUM_THRESHOLDは、もう1つの現在の部分集合を創るために
//残りのデータ部分集合内に十分なデータが残ることを保証する
//よう典型的に0.5である
/*****
或いは他にもし該ランダムな推定が該データ項目は現在のデータ部分集合へ行くべきと決
めたなら、NUMITEMSPERCATの望まれる割り当てが越えられたかどうかをチェックして見る
。もしそうでないなら、現在のデータ部分集合にデータ点を追加し、CountinStateをイン
クレメントする。
*****/
/*****
elseif (guess <= Threshold (output) ) {
    if (CountinState (output) < NUMITEMSPERCAT) {
        CopyRecordtoDataSubset;
        IncrementCountinState (output) ;
    }
    else {
        CopyRecordtoRemainingData;
        IncrementCountinRemainingState (output) ;
    }
}

/*****
又は最後に、もし該ランダムな推定が該データ項目が該残りのデータ部分集合内に行くべ
きことを決めるならば、該残りの部分集合用割り当てが越えられたかどうかをチェックす
る。もしそうでないなら、該残りのデータ部分集合へ該データ項目を追加する。もし該割
り当てが越えられたなら、もしそのカテゴリー内でより多くの項目が必要なら該データ項
目を該現在のデータ部分集合に追加する。
*****/

```

10

20

30

```

    elseif(CountinRemainingState (output) < (1-Threshold (output) ) *
        TotalCountinState (output) ) {
        CopyRecordtoRemainingDataSubset;
        IncrementCountinRemainingData (output) ;
    }
    elseif(CountinState (output) < NUMITEMSPERCAT){
        CopyRecordtoDataSubset;
        IncrementCountinDataSubset (output) ;
    }
} //データ記録ループの終了 10
//BalanceRandomizeの終了
NaturalRandomize
SampleSize=NumberOfDataRecords/NumberOfModels;
Threshold=1-SampleSize/NumberOfRemainingDataRecords;
Loop through output state {
    InitializeCountinState (output) to 0;
    InitializeCountinRemainingState (output) to 0;
}
Loop through data records {
    Loop through input features {
        if (input feature==min) {
            if (input FeatureMinFlag=CLEAR){
                IncludeTrainFlag=TRUE;
                FeatureMinFlag =SET;
            }
        }
        elseif (input feature==max) {
            if (input FeatureMaxFlag=CLEAR){
                IncludeTrainFlag=TRUE;
                FeatureMaxFlag =SET;
            }
        }
    }
    //フィーチャーループ終了
    output=ReadOutputState; //記録用に出力状態を読み出す
    guess=GuessRandomValue;
    / * * * * *
    もしデータ記録がフィーチャーの最小又は最大値の最初の場合なら、該データ部分集合及び残りのデータ部分集合の両者に記録をコピーする。
    / * * * * *
    if (IncludeTrainFlag=TRUE) {
        //該データ部分集合と該残り
        //のデータ集合との両者に記
        //録をコピーする
        CopyRecordtoCurrentDataSubset;
        CopyRecordtoRemainingDataSubset;
    }
    / * * * * *
    又はもし該ランダムな推定が該データ項目が該残りのデータ部分集合内に行くべきことを決めるなら、そのカテゴリー用に該残りの部分集合の統計的限界が越えられたかどうかをチェックする。もし越えられないならば、該残りのデータ部分集合に該データ項目を追加する。もし該割り当てが越えられたなら、該データ部分集合に該データ項目を追加する。 50

```

```

/ * * * * *

```

```

    elseif (guess <= Threshold) {
        if (CountinRemainingState (output) <
            Threshold*TotalCountinState (output) )
            CopyRecordtoRemainingDataSubject;
        else
            CopyRecordtoCurrentDataSubject;
    }

```

```

/ * * * * *

```

又はもし該ランダムな推定が該データ項目が現在のデータ部分集合内に入るべきことを決めるなら、そのカテゴリー用に該現在の部分集合の統計的限界が越えられたかどうかをチェックする。もしそうでないなら、該現在のデータ部分集合に該データ項目を追加する。もし該割り当てが越えられたなら、該残りのデータ部分集合に該データ項目を追加する。

10

```

/ * * * * *

```

```

    else{
        if (CountinState (output) <
            (1-Threshold)*TotalCountinState){
            CopyRecordtoCurrentDataSubject;
        }
        else
            CopyRecordtoRemainingDataSubject;
    }

```

20

```

} //データ記録ループ終了

```

／NaturalRandomizeの終了

EvolveFeatures

```

    SelectRandomStackofGenes(N);
    Loop Through each gene in Stack {

```

```

/ * * * * * 遺伝子から部分空間を創る * * * * * /

```

```

        ReadParameters();
        ReadSubspaceAxesfromGene();
        if (AdaptiveNumberofBinsFlag=SET)
            CalculateAdaptiveNumbins;
        else
            UseNumBinsinParameterList;
        if (AdaptiveBinPositionsFlag=SET)
            CalculateAdaptiveBinPositions;
        else
            CalculateFixedBinPositions;

```

30

```

/ * * * * * : 遺伝子から部分空間を創ることの終了 * * * * * /

```

```

        ProjectTrainDataintoSubspace;
        CalculateGlobalEntropyforSubspace;
    }
    EvolveGenesUsingGlobalEntropy();    // 遺伝子ループの終了
    // 遺伝的アルゴリズム
}
CreateTrainTestSubsets
    DetermineRangesofInputs;
    RandomizeTrainTestSubsets;
RandomizeTrainTestSubsets
{
    Threshold=ReadThresholdfromParameterList;
    Loop through data records in Data Subset {
        Loop through input features {
            if (input feature==min) {
                if (input FeatureMinFlag=CLEAR){
                    IncludeTrainFlag=TRUE;
                    FeatureMinFlag  =SET;
                }
            }
            else{
                if (input feature==max) {
                    if (input FeatureMaxFlag=CLEAR){
                        IncludeTrainFlag=TRUE;
                        FeatureMaxFlag  =SET;
                    }
                }
            }
        }
    }
    // フィーチャループの終了
    output=ReadOutputState;    // 記録用に出力状態を読み出す
    guess=GuessRandomValue;

```

10

20


```

if (guess<= Threshold){
    if (CountinTrainDataSubset (output) <
        Threshold (output) *TotalCountinState
        OR IncludeTrainFlag=TRUE)
        CopyRecordtoTrainDataSubset;
    else
        CopyRecordtoTestDataSubset;
}
else{
    if (CountinTestDataSubset (output) <
        (1-Threshold)*TotalCountinState (output)
        AND IncludeTrainFlag=FALSE){
        CopyRecordtoTestDataSubset;
    }
    else
        CopyRecordtoTrainDataSubset;
}
}
//データ記録ループの終了
//RandomizeTrainTestSubsetsの終了
ModelEvolution
{
    GenerateRandomStackofModelGenes(); //モデル遺伝子が遺伝子のク
                                        //ラスタであるランダムモ
                                        //デル遺伝子を発生させる
    Loop through each model gene in stack {
        CalculateMGFF(); //モデル遺伝子適応度関数
                        // {エムジーエフエフ (MGFF)}
                        //の計算
    }
    EvolveFittestModelGene(); //モデル遺伝子ループの終了
                            //最適モデル遺伝子を発展さ

```

10

20

```

        // せるため遺伝的アルゴリズム
        // をドライブするようエムジー
        // エフエフを使用
    }
    CalculateMGFF—モデル遺伝子適応度関数（エムジーエフエフ）の計算
    {
        IdentifyFeatureGenes();          // フィーチャー遺伝子の集合を
        // 同定するためモデル遺伝子を
        // パース (parse) する

        Loop through each feature gene {
            CreateFeatureSubspace();
            Loop through each test record {
                ProjectTestRecordintoSubspace();
                UpdateTestRecordPrediction();
            }
        }
        Total_Error=0;
        Loop through each test record {
            if (RecordPrediction!=ActualRecordOutput)
                TotalError=TotalError+1;    // インCREMENT誤差
        }
        MGFF=Total_Error;
    }

```

本発明の好ましい実施例がここで説明された。付属する請求項により規定された本発明の真の範囲から離れることなく変更や変型が該実施例内で行われ得ることは勿論理解されるべきである。本実施例は好ましくは、コンピュータで実行可能なソフトウェア命令のセットとしてソフトウェアモジュール内で説明された方法を実施するロジックを含むのがよい。中央処理ユニット（「シーピーユー（CPU）」）、又はマイクロプロセッサは該トランシーバーの動作を制御する該ロジックを実行する。該マイクロプロセッサは説明された機能を提供するために当業者によりプログラムされ得るソフトウェアを実行する。

【 0 1 9 0 】

該ソフトウェアは、磁気ディスク、光ディスク、そして該シーピーユーにより可読な何等の他の揮発性 [例えば、ランダムアクセスメモリー { 「ラム（RAM）」 }] 又は不揮発性 [例えば、読み出し専用メモリー { 「ロム（ROM）」 }] ファームウェア記憶システムを含むコンピュータ可読の媒体上に保持される 2 進のビットのシーケンスとして表され得る。データビットが保持される該メモリー配置も又該記憶されるデータビットに対応する特定の電氣的、磁氣的、光学式又は有機的特性を有する物理的配置を有している。ソフトウェア命令はメモリーシステムを有する該シーピーユーによりデータビットとして実行され、該電気信号表現の変換と該メモリーシステム内のメモリー位置でのデータビットの保持をもたらし、それにより該ユニットの動作を再構成させるか又は他の仕方に変えさせる。該実行可能なソフトウェアコードは、例えば、上記説明の様な方法を実施してもよい。

【 0 1 9 1 】

ここで説明されたプログラム、過程、方法そして装置は、他のように指示されてない限り、どんな特定の種類のコンピュータ又はネットワーク装置（ハードウェア又はソフトウェア）にも関係付けられず、限定されないことは理解されるべきである。種々の種類の汎用又は専用コンピュータ装置又は計算装置がここで説明された開示に依って使用されてもよく、動作を行ってもよい。

【 0 1 9 2 】

本発明の原理が適用される広範な種類の実施例を見ると、図解された実施例は単に例示的で本発明の範囲を限定すると取られるべきでないことを理解すべきである。例えば、本発

明は金融サービス市場、宣伝及びマーケティングサービス、製造過程に関連するシステム又は大きなデータ集合を有する他のシステムで使用されてもよい。加えて、該流れ線図の過程は説明されたものとは他のシーケンスで用いられてもよく、そして該ブロック線図ではより多く又はより少ない要素が使われてもよい。

【 0 1 9 3 】

ハードウェア実施例は種々の異なる形式を取ってもよいことは理解されるべきである。該ハードウェアはカスタムゲートアレー（custom gate array）または特定用途向け集積回路（application specific integrated circuit）{ “エイシック（ASIC）” } で集積回路として実施されてもよい。勿論、該実施例は個別ハードウェア部品（discrete hardware components）と回路で実施されてもよい。特に、ここに説明した論理構造と方法の過程はエイシックの様な専用ハードウェアで、又はマイクロプロセッサ又は他の計算素子により行われるプログラム命令として実施されてもよい。

10

【 0 1 9 4 】

請求項はその効果に対し述べられていない限り要素の説明された順序に限定されるとして読まれるべきでない。加えて、何れの請求項でも用語 “手段（means）” の使用は 3 5 ユー・エス・シー・§ 1 1 2、パラグラフ 6 を行使するよう意図されており、該用語 “手段” を有しない何れの請求項もそのように意図されてない。従って、下記請求項の範囲と精神に入る全ての実施例とその等価物は本発明として請求されている。

【図面の簡単な説明】

【図 1】 本方法の全体的流れを図解するブロック図である。

20

【図 2 A 及び 2 B】 適合型ビンニングの例を示す。

【図 2 C】 データバランシングの方法を示す。

【図 3 A】 1 次元のフィーチャー部分空間を示す。

【図 3 B】 2 次元のフィーチャー部分空間を示す。

【図 3 C】 3 次元のフィーチャー部分空間を示す。

【図 4】 どの入力フィーチャー部分空間に含まれるかを表す例示的 2 進ビット記号列を示す。

【図 5 A 及び 5 B】 “情報豊富な” 入力フィーチャーの発展を図解するブロック線図である。

【図 5 C】 2 進記号列適応度の加重ルーレット選択ホイール（weighted roulette wheel）を示す。

30

【図 5 D】 交叉（crossover）操作線図を示す。

【図 6】 ローカルエントロピーパラメータを計算する方法を図解するブロック線図である。

【図 7】 グローバルエントロピーパラメータを計算する方法を図解するブロック線図である。

【図 8】 ローカル及びグローバル情報コンテンツの計算を図解する。

【図 9】 ローカルエントロピーパラメータとグローバルエントロピーパラメータの例を示す。

【図 1 0 A】 最適モデルを決定する方法を図解するブロック線図である。

40

【図 1 0 B】 モデル発展の方法を図解するブロック線図である。

【図 1 1】 情報写像（information map）を発生させる方法を図解する。

【図 1 2】 遺伝子リストとそれの付随情報写像の例である。

【図 1 3】 エグゾースチブな次元のモデリング過程の方法を図解するブロック線図である。

【図 1 4】 出力状態確率ベクトル / 出力状態値を計算する過程の方法を図解するブロック線図である。

【図 1 5】 モデル遺伝子用適応度関数を計算する方法を図解するブロック線図である。

【図 1 6】 1 つのフレームワークを発展させるために分布状階層的モデリングの方法を図解するブロック線図である。

50

【図 17 A 及び 17 B】 フレームワーク発展の方法を図解するブロック線図を含む。

【図 18 A】 スーパーフレームワークを発展させるための分布状モデリングの方法を図解するブロック線図である。

【図 18 B】 スーパーフレームワーク発展用の考慮点のリストである。

【図 19 A 及び 19 B】 クラスタ発展の方法を図解するブロック線図である。

【図 19 C】 データクラスターを発見する方法を図解するブロック線図である。

【図 19 D】 画像的表現用グローバルクラスタリング指数の計算方法を図解するブロック線図である。

【図 1】

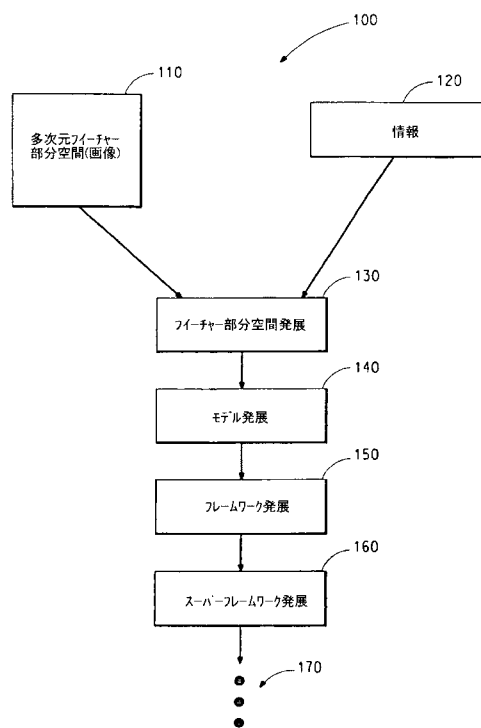


FIG. 1

【図 2 A - B】

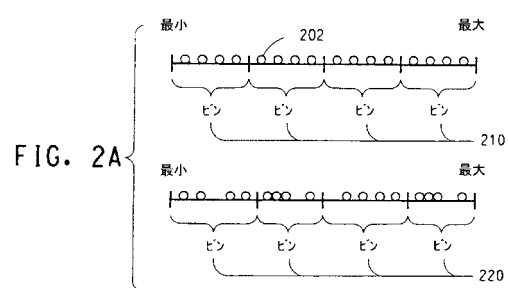


FIG. 2A

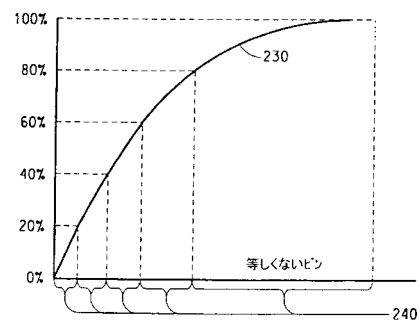


FIG. 2B

【図 2 C】

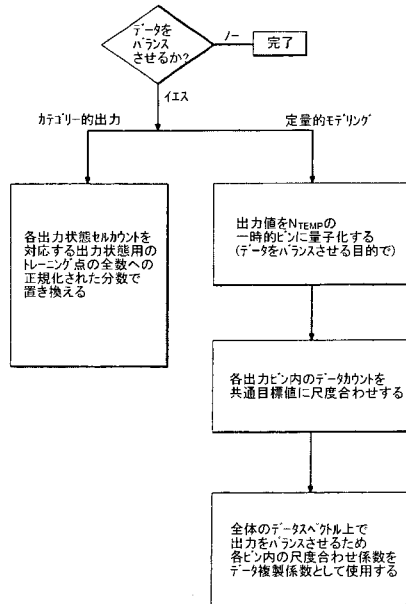


FIG. 2C

【図 3 A - B】

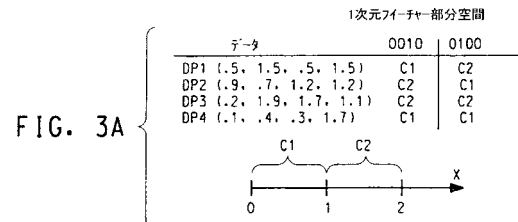
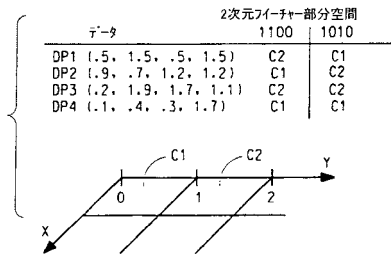
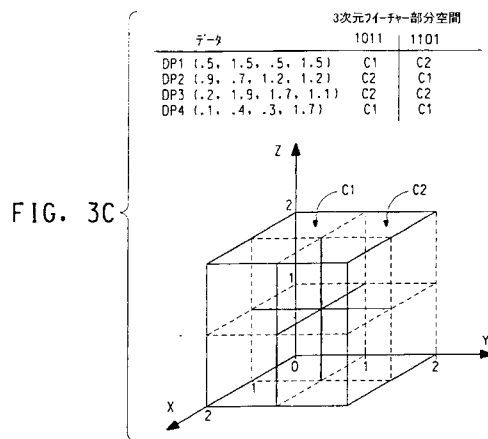


FIG. 3B



【図 3 C】



【図 4】

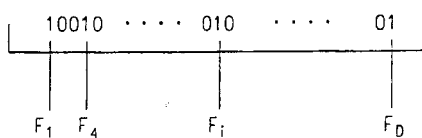
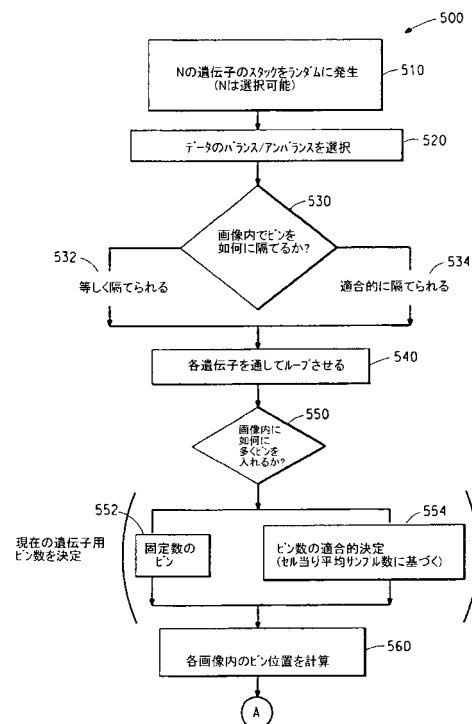


FIG. 4

【図 5 A】



【図 5 B - D】

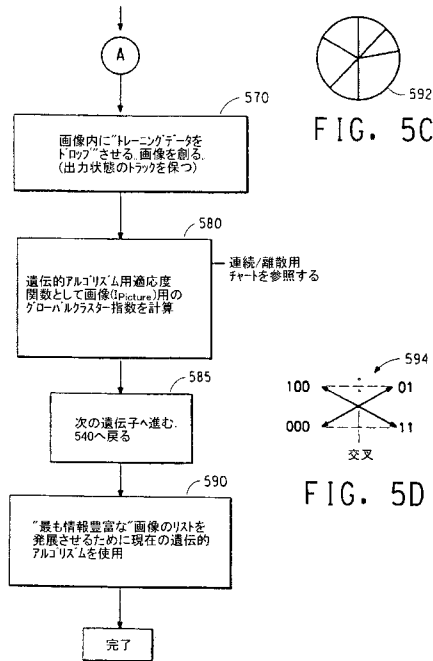


FIG. 5B

【図 6】

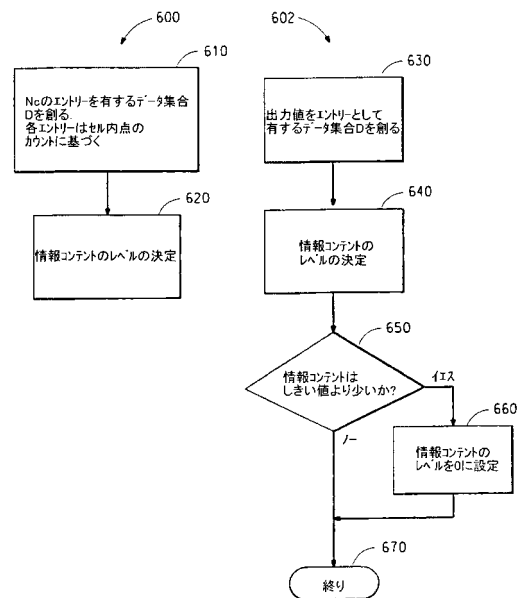


FIG. 6

【図 7】

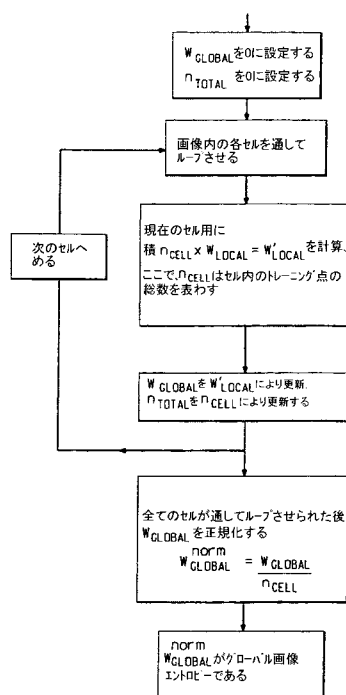


FIG. 7

【図 8】

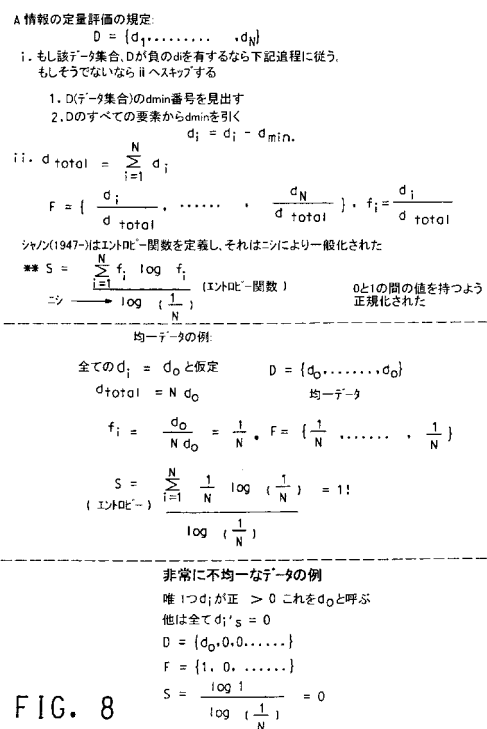
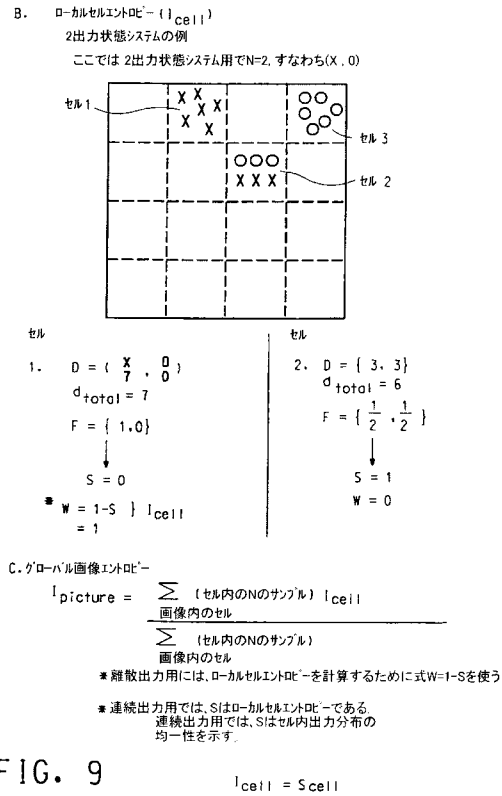


FIG. 8

【図 9】



【図 10 A】

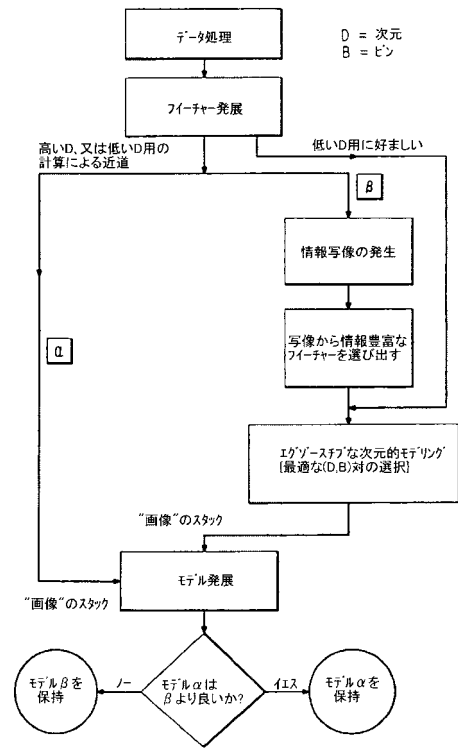


FIG. 10A

【図 10 B】

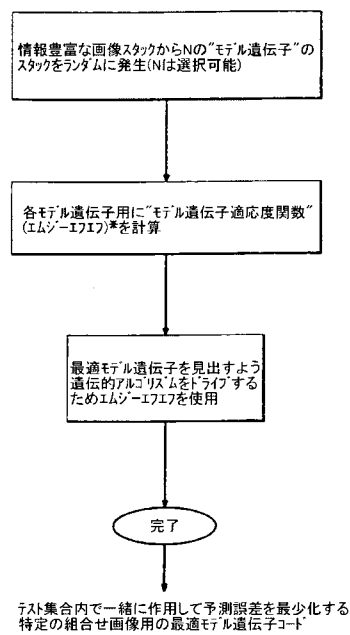


FIG. 10B

【図 11】

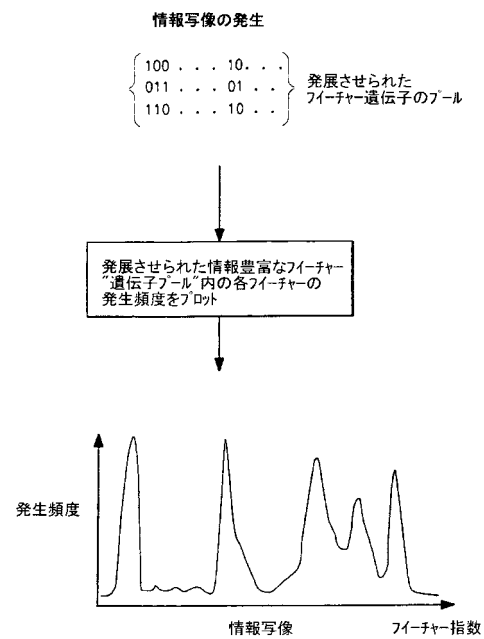
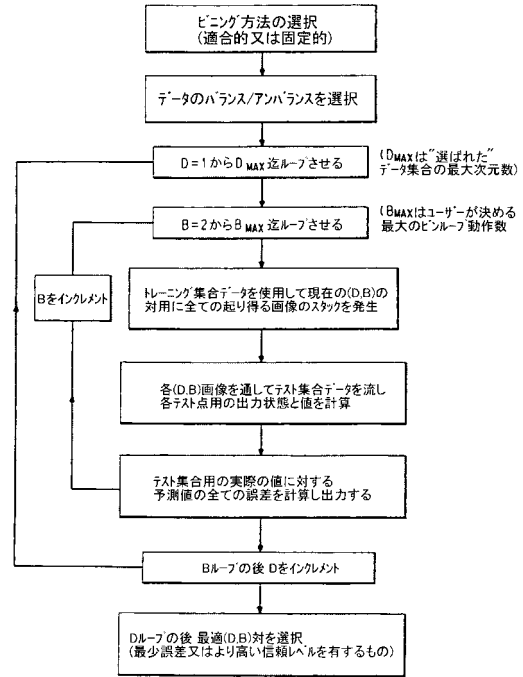


FIG. 11

【図 12】



【図 13】



【図 14】

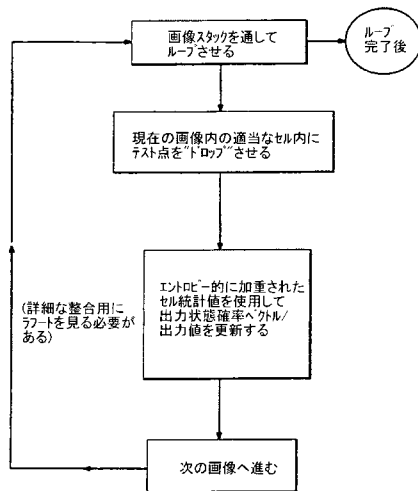


FIG. 14

【図 15】

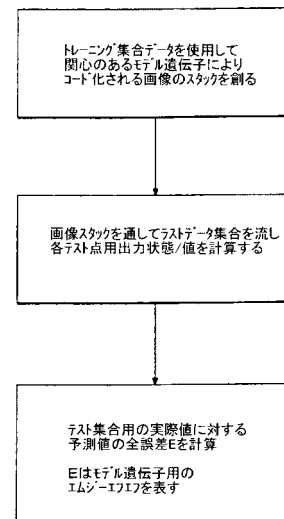


FIG. 15

【図 16】

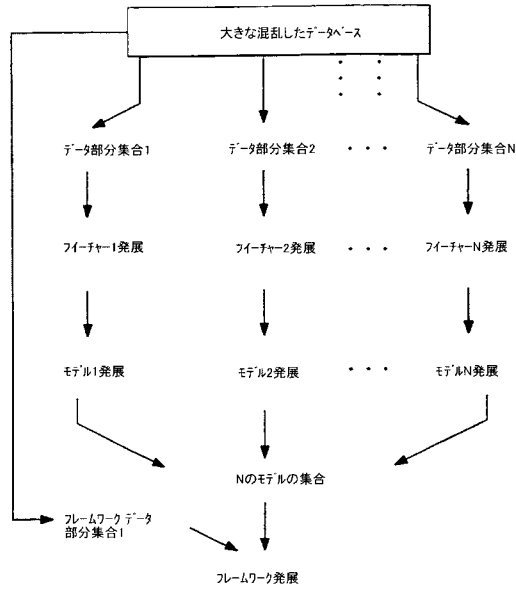


FIG. 16

【図 17 A】

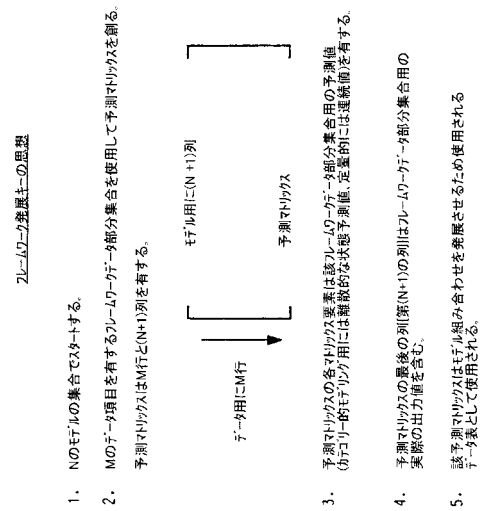


FIG. 17A

【図 17 B】

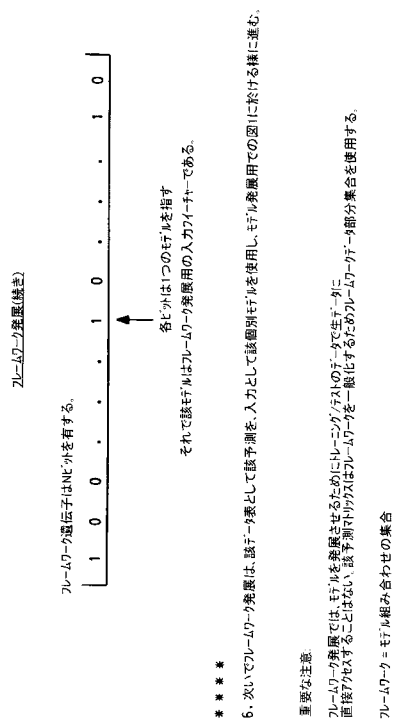


FIG. 17B

【図 18 A】

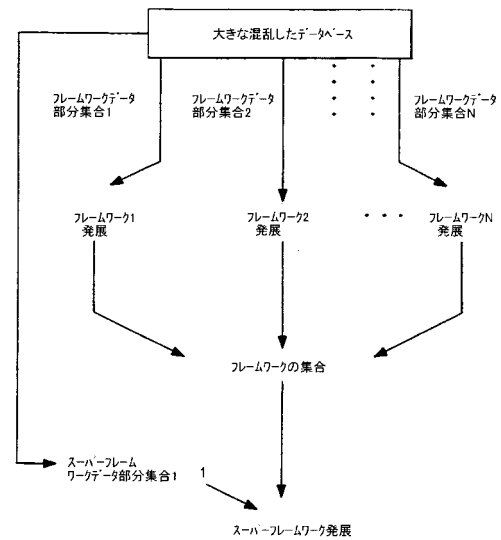
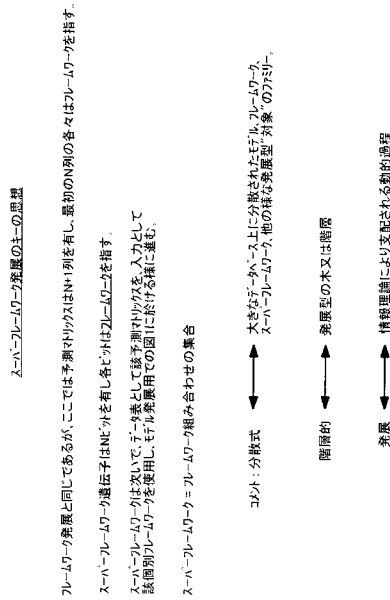


FIG. 18A

【図 18 B】



発展は、得られるべき情報がこれ以上ないか又は我々が望む終了点に到着するまで、該階層の下へ進む。

FIG. 18B

【図 19 A】

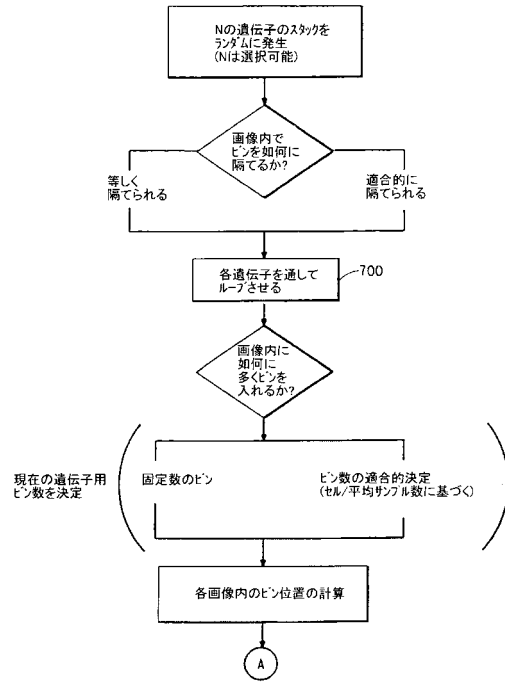


FIG. 19A

【図 19 B】

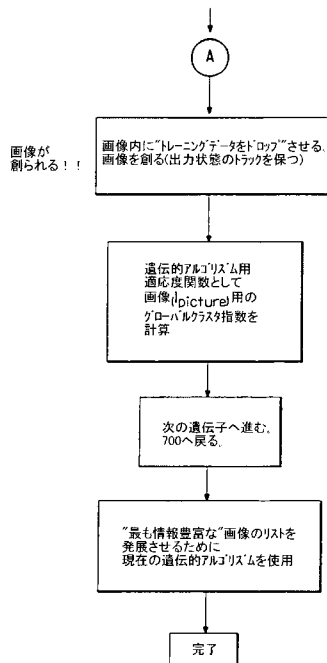


FIG. 19B

【図 19 C】

基本的動機:

インフォホルブフレームワークは入力フィーチャーの類似性に基づきデータクラスターを発見するため使い得る。これはどんな出力状態からも独立しており、入力フィーチャーの類似性に基づき大きなデータベースのより小さい部分集合への部分分割に非常に有用である。

最高レベルの流れ図

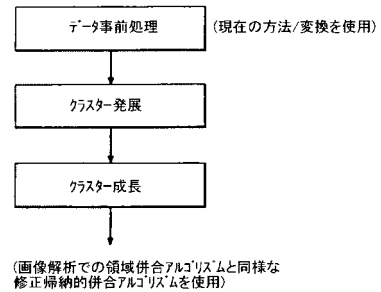


FIG. 19C

【図 19D】

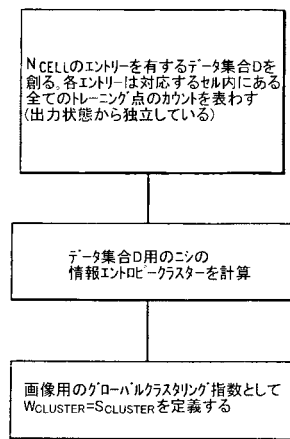


FIG. 19D

フロントページの続き

- (72)発明者 バイドヤナサン, アクヒルスウオー・ガネシユ
アメリカ合衆国デラウェア州 1 9 7 0 7 ホツケシン・ロビンコート 4 4
- (72)発明者 オーエンス, アーロン・ジエイ
アメリカ合衆国デラウェア州 1 9 7 1 3 ニューアーク・ルネイブレイン 2 3 ・シルバーブルツク
- (72)発明者 ウイトコム, ジエイムズ・アーサー
アメリカ合衆国ノースカロライナ州 2 8 7 1 2 ブルバード・カントリークラブロード 1 3 1 5

審査官 長谷川 篤男

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

- (58)調査した分野(Int.Cl. , D B 名)

G06N 3/00

IEEE Xplore

JSTPlus(JDreamII)