

高次元最適ピラミッドを用いた数値属性ルールの生成とデータマイニングへの応用

Higher-Dimensional Pyramid Construction Problem and Application to Data Mining

全 眞嬉¹ Danny Z. Chen²

加藤 直樹³ 徳山 豪⁴

Jinhee CHUN

Danny Z. CHEN

Naoki KATO

Takeshi TOKUYAMA

福田らによって提案された最適領域ルールは数値属性データで構成されるデータベースに対する有効なデータマイニングツールである[FMMT99, FMMT01]。しかしながら、従前の方法では2つの欠点がある:(1)各ルールは高々2変数の数値属性に対応でき(2)判定は与えられたデータの正確な位置ではなく、単純に領域Rの中にあるか外にあるかだけに基き行われる。本論文では、これらの欠点を取り除くための新しい方法の提案をする。具体的にはグラフアルゴリズムを用いて、2つ以上の属性を持つ最適数値属性結合ルールと階層構造の数値属性結合ルールを与える。又、本論文の手法は異常なデータの除去とデータクラスタリングに適切である。

Optimized region rules developed by Fukuda et al.[FMMT99, FMMT01] are effective tools for data mining in databases with numeric data. However, there are two drawbacks in the previous methods: (1) each rule can contain at most two numeric conditional attributes, and (2) the decision is made based only on whether a given data is inside or outside a region R, but not on the exact position of the data. In this paper, we propose a new method for removing these drawbacks. Indeed, by applying graph algorithms, we give optimized numeric association rules with more than two attributes, and give layered-structure numeric association rules. Our method is also applicable to removal of exceptional data and data clustering.

1. はじめに

データマイニングの手法として「AならばBである」といったデータベースの属性間の相関関係を求める結合ルールがあり、Aprioriアルゴリズム[AIM]が提案された。現行のデータベースには性別、血液型といったカテゴリ属性と年齢、体重といった数値でとられる数値属性が含まれているが、Aprioriアルゴリズムはカテゴリ属性間の結合ルールを求めるためのものであり、数値属性に対しては直接適用することはできない。なぜなら、数値属性は値に順序があるため、値を連続した区間で表すことが重要であるが、Aprioriアルゴリズムではそのようなこと

1 学生会員 東北大学大学院情報科学研究科システム情報科学専攻 jinhee@dais.is.tohoku.ac.jp

2 非会員 University of Notre Dame, Dept. of Computer Science and Engineering chen@cse.nd.edu

3 非会員 京都大学大学院工学研究科建築工学専攻 naoki@archi.kyoto-u.ac.jp

4 非会員 東北大学大学院情報科学研究科システム情報科学専攻 tokuyama@dais.is.tohoku.ac.jp

は考慮していないからである。よって、数値属性を含む結合ルールを求めるには別のアルゴリズムを用いる必要がある。

福田ら[FMMT01]は領域切り分けを用いて数値属性とカテゴリ属性間の結合ルールの生成を行う手法を提案している。この手法では「 $x \in R$ ならば B である」(R はいくつかの数値属性の値の空間の部分領域)のようなルールを発見している。

福田らによる領域として切り出しを行うと、領域に入っていないかいないかだけで判定を行うため、境界線上のデータと中央部のデータが同じ扱いをされるといった位置情報の損失がある。領域切り分けを用いて生成された結合ルールを学習に用いると、制限が強く、学習データに入っているノイズ(離れ値)を排除してしまうため、学習データ以外のデータに対して生成されるルールは不自然な領域ルールを生成するといった、過学習の問題点がある。また、属性を3つ以上与えると、計算量が多く計算が困難になる。すなわち、属性次元の制限の問題がある。ここで、本来は切り取るべきものは多次元正規分布のような特徴をもったデータ分布であり、領域として切り取るより、性質の良い関数として捉える事が望ましい。

本論文では、以上の問題を解決するために、最適ピラミッドによる結合ルールの表現法を提案する。つまり、領域として切り出すのではなく、性質の良い関数として切り出しを行う。数値属性または複数の数値属性についての最適な階層構造でデータの近似を行い、階層的な数値結合ルール、すなわち最適ピラミッドを用いた数値属性結合ルールの提案をする。

手段として、幾何学的なアルゴリズムを用いた最適ピラミッド構築問題への定式化を行う。最適ピラミッドとは、入力データが持っている位置情報の損失を最小にする(最小の近似エラーを持つ)、すなわち、パラメトリックゲインを最大にする領域を積み重ねた単峰な図形への変更である。

2. 結合ルール

結合ルールは、タプル内の属性間の相関関係である。X, Y が属性に関する命題式のとき、 $(X \rightarrow Y)$ を結合ルールと呼ぶ。「X ならば Y である」という意味である。Y は通常、単一属性に関する命題式である。

健康診断データベースを例として考える。正常体重域を超えた体重を肥満、正常血糖値を超えた血糖値を異常血糖値と呼ぶことにする。「肥満で異常血糖値である人は糖尿病の検査が必要である」という診断法則があるとする。健康診断データベースの診断法則の結合ルールは(肥満 = yes \wedge 異常血糖値 = yes) \rightarrow (糖尿病検査 = yes) と表現される。

X が単一属性に関する命題式のとき、 $(X \rightarrow Y)$ は1次元結合ルールと呼び、 $(X_1 \wedge X_2 \rightarrow Y)$ のように前提条件に X_1, X_2 の2つの属性を用いた場合は2次元結合ルールと呼ぶ。

結合ルールの尺度である、確信度とサポートについて述べる。結合ルール $X \rightarrow Y$ の確信度(confidence)とは、データベース D の X を含むトランザクションのうち、Y を同時に含む割合のことである。すなわち、 $|X \wedge Y|/|X|$ が c%であることを、 $\text{conf}(X \rightarrow Y) = c\%$ と表記する。「結合ルール $X \rightarrow Y$ は c%の確信度である」という意味を持つ。結合ルール $X \rightarrow Y$ のサポート(support)とは、データベース D の(X)を含むトランザクションの、トランザクション全体に対する割合のことである。

このとき、 $\text{support}(X \rightarrow Y) = s\%$ と表記する。「結合ルール $X \rightarrow Y$ は s%のサポートである」という意味を持つ。

各項目間の関連に説得力のある根拠を提示することで、抽出された規則はユーザに判断基準を与える説明性を持つ。このような規則の根拠を確信度とサポートを用いて表す。

本研究では、新しい領域族を階層最適化し、より高次元のルールの効率的な生成を行う。本研究の結果は、結合ルール生成だけでなくデータビジュアル化およびデータマイニングへのファジーアプローチにおいても有効に応用できる。階層構造は決定論的な決定ルール($x \in R \rightarrow C = \text{yes}$)に比較して、強いルールの影響を縮小する方法を適用する。拘束力の弱いルールで判定をする、非決定性を持たせた柔軟な決定システムの構築を行う。

3. 最適ピラミッド問題

階層的結合ルール生成問題を定式化するためにより一般の幾何学問題を与え、データマイニングへの利用をその応用として行う。

μ と ρ を $n = m^d$ セルの d 次元ボクセルグリッド Γ 上の非負整数値関数とする。すべてのセル $c \in \Gamma$ に対し $\mu(c) \geq n$ かつ $\mu(c) \geq \rho(c)$ と仮定する。データマイニングへの応用では μ はセル c に対応するデータの数、すなわちサポート(support)の関数を表し、 ρ はセル c において条件属性と目的属性を同時に満たすデータの数、すなわちヒット(hit)関数を表す。また、確信度(confidence)はサポートをヒットで割ったもので、 $\text{conf}(c) = \rho(c) / \mu(c)$ とする。

N に領域族 F を固定する。一般性を失わず、 $\emptyset \in F$ および $\Gamma \in F$ と仮定できる。

Definition 1 $P(t_0) = \Gamma$ および $t > t'$ のとき $P(t) \subseteq P(t')$ を満たす F の領域の列 $P = P(t_i)$ ($i=0,1,2, \dots, h$) を考える。ここで $t_0 < t_1 < t_2 < \dots < t_h$ は高さと呼ばれる実数である。 $P_{t_{i+1}} = \emptyset$ で $\text{conf}(P(t_i) \cdot P(t_{i+1})) = t_i$ の場合、 P を μ に関して $\text{conf} = \rho / \mu$ を近似する。 ρ に近似するピラミッド(あるいはピラミッド構造)と言う。 P の近似エラーは

$$\sum_{i=0}^h \sum_{c \in P(t_i) - P(t_{i+1})} (\text{conf}(c) - t_i)^2 \mu(c)$$

によって定義される。これは確率密度関数として μ を考えたとき、確信度 ρ / μ と、 $[f_P(x) = t_i \text{ if } x \in P(t_i) \cdot P(t_{i+1})]$ によって定義されたピラミッドの表面関数 f_P の間の 2 乗された L_2 距離である。すべてのピラミッド中で最小近似エラーを持つとき、ピラミッド P を最適であると言う。

μ がユークリッド空間の体質を与える密度関数であるとき、最適ピラミッドは、位置ポテンシャルの損失を最小限にする ρ / μ の単峰な変更と見なすことができる。これは計算幾何学および地理学(特に $d=2$ の場合)の基礎的な問題となる。最適ピラミッドの構築は領域切り出し問題の自然な拡張であり、データマイニングに加えて多くのアプリケーションにおいて有用になる。

3.1 最適ピラミッドとパラメトリックゲイン

Lemma 1 領域を値に持つ写像 $P(t): (0, \infty) \rightarrow F$ で、単調条件 $[t > t' \text{ のとき } P(t) \subseteq P(t')]$ を満たし、目的関数 $J(P) = \int_0^\infty (\rho(P(t)) - t \cdot \mu(P(t))) dt$ を最大にする関数 $P(t)$ を考える。□

このとき、 $P(t)$ は $t \in (t_{i-1}, t_i)$ のとき $P(t) = P(t_i)$ 及び $t > t_h$ のとき $P(t) = \emptyset$ を満たす、高々 $n+1$ の変化値 $0 = t_{-1} < t_0 < t_1 < \dots < t_h$ を持ち、さらに $P(t_0), P(t_1), \dots, P(t_h)$ から成る P は最適ピラミッドである。

Proof Γ に $n = m^d$ 個のピクセルがあるので、 $t \in n(0, \infty)$ のとき $P(t)$ の中に高々 $n+1$ の異なる領域がある。したがって、変化値の数が $n+1$ 以下であることは明らかである。目的関数を最大にするために、 $\rho(P(t_i)) - t \cdot \mu(P(t_i))$ と $\mu(P(t_i)) = \rho(P(t_{i-1})) - t \cdot \mu(P(t_{i-1}))$ は変化値 $t = t_{i-1}$ で等しい $t_{i-1} = \text{conf}(P(t_{i-1}) \setminus P(t_i))$ であり、 P は ρ のピラミッドである。2 乗された L_2 エラー $E(P) :=$

$$\sum_{i=0}^h \sum_{c \in P(t_i) - P(t_{i+1})} (\text{conf}(c) - t_i)^2 \mu(c)$$

を最大にする関数 $P(t)$ は、 $E(P)$ を最小にするピラミッド P を定義する。よって証明された。

したがって、直観的に最適ピラミッド P は、できるだけ大きなパラメトリックゲインを持つ横断面を積み重ねることにより得られる。パラメトリックゲイン $g_t(R, \rho, \mu)$ を最大にする F の領域 $R^{\text{opt}}(t)$ を考える。

直感的には t が増加する場合、 $R^{\text{opt}}(t)$ は減少する。 $\{R^{\text{opt}}(t)\}$ (正確には、関数 $R^{\text{opt}}(t)$ の閉包)がピラミッドを形成する場合、それは明らかに最適ピラミッドである。ところが、最大ゲイン領域 $R^{\text{opt}}(t)$ を積み重ねるとき、 $R^{\text{opt}}(t) \subset R^{\text{opt}}(t')$ が $t > t'$ に対して成立するとは限らないので、それらは必ずしもピラミッドを形成するとは限られない。これは、最適ピラミッドを計算することを非常に難しくする。

3.2 閉集合族に対する考察

ある集合族 S が $G = [0, n]^d$ の閉集合族であるとは、 $G \in S, \emptyset \in S, [X, Y \in S \rightarrow X \cap Y \in S]$ かつ $[X, Y \in S \rightarrow X \cup Y \in S]$ である事を言う。

Theorem F が閉集合族であれば、各々の t に対し $\rho(R) - t \cdot \mu(R)$ を最大にする領域 $R \in F$ を $\Psi(t)$ とすると、 Ψ は単調条件を満たし、従って、最適ピラミッドになる。□

Proof $t > t'$ のとき $A = \Psi(t) \subseteq B = \Psi(t')$ である。 $R \subset R'$ の場合、 t に対して関数 $g_t(R, \rho, \mu) - g_t(R', \rho, \mu)$ は増加しない関数である。 $A \setminus B$ が空集合でない場合、 $0 \geq g_{t'}(A \cup B, \rho, \mu) - g_{t'}(B, \rho, \mu) = g_{t'}(A, \rho, \mu) - g_{t'}(A \cap B, \rho, \mu) \geq g_t(A, \rho, \mu) - g_t(A \cap B, \rho, \mu)$ である。 $A = \Psi(t)$ なので、 $g_t(A, \rho, \mu) \geq g_t(A \cap B, \rho, \mu)$ になる。従って $0 \geq g_t(A, \rho, \mu) - g_t(A \cap B, \rho, \mu) \geq 0$ であり、 $g_t(A, \rho, \mu) = g_t(A \cap B, \rho, \mu)$ になる。従って $A = A \cap B$ であり、 $A(B)$ である。□

従って、集合族 F が閉集合族で構築される場合、 F の最適のピラミッドを計算するための効率的なアルゴリズムの設計が必要である。

4. 2次元の場合

n をピクセルの数、 N をサポート値の和(すなわち、データの数)とする。2次元の領域族でもっとも基本的なのは軸方向長方形の族の場合である。しかし、この場合では $O(n^5)$ の計算時間のアルゴリズムしか知られていない[CKT02]。ここでは、実用的な多くの領域族についてより効率的な手法を考える。

Definition 2 2次元ピクセル面 $G = [0, m]^2$ ($m = \sqrt{n}$) 内で、適当な関数 $y = h(x)$ より真に下にあるピクセル全体の和になっている領域を下半切断領域と呼ぶ。□

Definition 3 2次元ピクセル面 $G = [0, m]^2$ 内の1点 p に対し、 p を含む長方形たちの和集合を p を中心にする矩形和楕円型領域 (rectilinear ellipsoid) と呼ぶ。□

点 p を中心とする矩形和楕円型領域の族は、点 p を含む矩形全体の族の集合和、集合積に対する閉包になっている。矩形和楕円型領域の族は直交凸領域の族の部分族であるが、軸方向楕円の離散化を含む広い領域族であり、指数個の領域を持つ。例として、単調減少関数の下半切断領域は原点を中心とする矩形和楕円型領域となる。

Lemma 2 点 p を中心とする矩形和楕円型領域全体の集合は閉集合族である。□

Proof 矩形和楕円型領域は点 p を含む長方形の和集合としてとられる領域である。点 p を中心とする矩形和楕円型領域 F とし、 F の任意の領域 R と R' に対して、 $R \cap R' \in F$ および $R \cup R' \in F$ である。すなわち、 F の任意の領域 R と R' は和集合と共通集合に閉じているので、 R^d の領域族 F は閉集合族である。□

Lemma 3 与えられた点 p を中心とする矩形和楕円型領域で与えられた t に対して利得を最大にする領域は $O(n)$ で計算できる。□

Proof まず、 $p=(0,0)$ の場合を考える。この場合、矩形和楕円型領域は単調非減少関数で区切られた階段状の領域になる。まず、最初の行における最初の i 個のピクセルのゲインの和

(Prefix和)である $f(i, j) = \sum_{s=1}^i \rho((s, j)) - t \cdot \mu((s, j))$ を計算する。これはすべての (i, j) に対して $O(n)$ 時間で計算できる。 □

(i, j) を p が中心で (i, j) を含み、 y 座標が j より大きい領域を含まないような矩形和楕円型領域の中でゲインが最大になるもののゲインと定義する。すると、 p を中心とする矩形和楕円型領域の中でゲインが最大になるもののゲインは $\max_{j=1}^m (1, j)$ によって求まる。 (i, j) は以下のダイナミックプログラミングによりすべての (i, j) に対して $O(n)$ 時間で求まる:

$$(i, j) = (i-1, j) + f(i, j) \quad (1)$$

$$(i, j) = \max\{(i, j), (i, j+1)\} \quad (2)$$

そして、領域はダイナミックプログラミングプロセスのバックトラッキングにより計算できる。

一般の $p=(x_p, y_p)$ の場合は、ピクセル領域を $x=x_p$ と $y=y_p$ で4分し、各象限で上記と同じ操作を行い、得られた領域の和領域を求めればよい。

Theorem 2 指定された一点 p を中心とする矩形和楕円型領域全体の集合に関する最適ピラミッドは、 $O(n \log N)$ 時間で計算される。ここで N は問題の出力精度であり、 $1/N$ の近似精度で最大な積分値を持つ解が出力される。(従って、 N が入力精度以上なら最適解を出力する)。 □

Proof 閉集合族なので、各 t に関して、利得を最大にする領域を求めればよい。これは $O(n)$ 時間でできる。更に、ある t での最適領域を計算した時に、これが G を領域の内部と外部に分割するため、各部分を独立に、対応する部分に問題を縮小して解ける。従って分割統治が可能である。 t はプロセスの分岐を用いる事により、深さ $\log N$ の再帰で全ての部分での最適領域の計算が終了する。各深さでの全ての部分領域での最適領域は計算時間は $O(n)$ で求まる。

従って、全体に関する最適ピラミッドは深さが $\log N$ なので $O(n \log N)$ 時間で計算できる。 □

5. 高次元の場合

5.1 小さい領域の場合

F が M 個の異なる領域を持っている場合、最適ピラミッドは d 次元の場合 M と n についての多項式時間で計算できる。頂点集合と F とする閉路を持たない有向グラフ $H(F) = (V, E)$ を構築する。各 F のペア R と R' が $R \setminus R'$ を満たすとき、グラフに有向辺 $e = (R, R')$ を加える。そして $(R \setminus R') = t(e) \cdot \mu(R \setminus R')$ を計算する。値 $t(e)$ は e の高さラベルと呼び $r(e) = t^2(e)$ ($R \setminus R'$) / 2 は e の利得と呼ぶ。

$t(e_{i-1}) < t(e_i)$ が $i=1, 2, \dots, q$ で成立する場合、有向路 $p = e_0, e_1, \dots, e_q$ を許容有向路(admissible)と呼ぶ。許容有向路の利得は辺の利得の総和である。

Lemma 4 最適ピラミッドは、 (R, R') が経路上の辺である場合のみ $R \setminus R'$ がピラミッドの面である場合、 $H(F)$ の最大利得の許容路 (admissible path) に対応する。 □

したがって、最適ピラミッド問題を循環路を持たない有向グラフ $H(F)$ の最大重み経路(maximum-weight-path)問題にすることができる。 $H(F)$ の各有向閉路が大部分で n 辺を持っていることに注意する。ダイナミックプログラミングアルゴリズムを使うことによって、次の結果を得る:

Theorem 3 M 個の異なる領域の F のための最適ピラミッドは $O(M^2 n)$ 時間で計算できる。 □

しかし、上記のアルゴリズムは実用的ではない。例えば、長方形の領域族は $O(n^2)$ 領域を持っている。従って、上記の計算

時間は $O(n^5)$ である。さらに、正確な階層状の領域ルール計算のために、 M が n において幾何級数的に大きい族を考える。

したがって、特別な領域族のために、より効率的なアルゴリズムを考える。

5.2 直交領域の stabbed union

領域の典型的な閉集合族について考える。の固定セル c については、各々が含んでいる直交領域の結合として R を表わすことができる。の領域 R は c で直交領域の stabbed union と呼ぶ。セル c は R の中心セルと呼ぶ。2次元の場合は矩形和楕円型領域とほぼ同一で、相異はグリッド点の代わりにセルを考える点である。

セル c のすべての stabbed union 族が閉集合族であることは明らかである; 実際、 c を含んでいるすべての長方形集合族に対して閉じている。与えられたピラミッドは、中心セルでの stabbed union 族に基づいている。当然、中心セル(あるいは点)は、ピラミッドのピークである。最適ピラミッドを計算するためのアルゴリズムを設計するためには、stabbed unions 族の最大パラメトリックゲイン領域計算のために効率的なアルゴリズムを必要とする。上記の目的のために、問題を一般化し、グラフアルゴリズムを適用する。

5.3 グリッドグラフ推移は閉包

ボクセルグリッドにおいて、中心セル $c = (c_1, c_2, \dots, c_d)$ を固定し、頂点セットが のすべてのボクセルから成る有向グラフ $G(c)$ を定義する。ボクセル $p = (p_1, p_2, \dots, p_d)$ および $q = (q_1, q_2, \dots, q_d)$ について、 p と q の間の L_1 距離は $dist(p, q) = \sum_{i=1}^d |p_i - q_i|$ である。 p の近隣のセルとは p から L_1 距離が 1 のセルである。セル p およびその近隣 q について、有向辺が定義される。その方向は $dist(p, c) = dist(q, c) + 1$ の場合 (p, q) (つまり p から q まで) であり、そうでなければ (q, p) である。グラフ $G(c)$ は $d(m-1) \times m^{d-1} = O(n)$ 辺 (d は定数とする) の弱い連結有向グラフであり、また、 c はそのユニークな sink 頂点(つまり出発する辺のない頂点)である。

$G(c)$ の部分グラフ $H=(V, E)$ において、 V の各頂点 v から c までの H に有向路が存在する場合、 H を根付き部分グラフ(rooted subgraph)と呼ぶ。

$G(c)$ の根付き部分グラフ $H=(V, E)$ が与えられ、 H に v から u までの有向路が存在する場合、頂点 u が頂点 v によって H において支配されると言う。頂点の集合 W において、 W のすべての頂点 V に対し、 V に支配される頂点がすべて W に含まれるとき、 W を H -閉包と言う。各 H -閉包は c 中でセル c を含んでいる連結領域を定義する。

与えられた $G(c)$ の根付き部分グラフ H の、すべての H -閉包の集合族 F_H を考える。推移閉包であるという性質は集合和と集合積について閉じているので、次の命題が示せる:

Proposition 1 $G(c)$ の根付き部分グラフ H について、 F_H は閉集合領域族である。 □

次の補題は、 $G(c)$ および F_H の定義からなる。

Lemma 5 の中の領域 R は、 $G(c)$ -推移閉包の場合のみ c の stabbed union である。 □

Lemma 6 H と H' を $G(c)$ の全域根付き部分グラフとする。このとき、 H が H' の部分グラフならば $F_H \subseteq F_{H'}$ である。 □

5.4 F_H の最適ピラミッド計算アルゴリズム

$G(c)$ の根付き部分グラフ H を固定する。 F_H に関しての最適ピラミッドの面の高さを定義するパラメーター値 t を考える。また、各ボクセル p (H の対応する頂点) に重さ $(p) - t \cdot \mu(p)$ を与える。次の補題により、面の高さを与える t が小さな分母お

よび分子を持つ有理数であると仮定できる。

Lemma 7 t がある領域族のための最適ピラミッド面を定義する高さである場合, t は N 以下の 2 つの整数の比によって表わされる有理数である。□

Proof t が面 P_i を定義すると仮定すると, $(P_i \setminus P_{i+1}) = t \cdot \mu(P_i \setminus P_{i+1})$ である。と μ は N 以下の整数値をとるので, 補題は成り立つ。□

定義によって, 最大(パラメトリック)ゲイン領域 $P^{opt}(t)$ F_H は, 領域のボクセル重みの合計を最大限にする H -閉包である。グラフ理論的な用語では, 重み付き有向グラフ H の最大推移閉包であり, 辺の集合(cut set)の削除により c を含んでいる H の連結成分として得られる。

Theorem 4 (Hochbaum)実数頂点重みを持つ n 頂点と m 辺からなる有向グラフ G の最大重み閉包は $O(T(n,m))$ 時間で計算できる。ここで, $T(n,m)$ は非負辺重みを持つ n 頂点と m 辺からなる有向グラフ H の最小 s - t カット計算時間である。□

Theorem 5 集合族 F_H のための最適ピラミッドは $O(n^{1.5} \log n \log^2 N)$ 時間で計算できる。□

Proof N 以下の N 整数の対によって定義された有理数の集合 S の中のすべての可能な面の高さを知るために, 2 分探索の一般化を適用する。プロセスは多くの部分プロセスに分岐するので, 2 分分岐探索と呼ぶ。 $t_0 < t_1$ に対し, $P^{opt}(t_0)$ と $P^{opt}(t_1)$ をすでに計算していると仮定する。

$P^{opt}(t_0) = P^{opt}(t_1)$ の場合, 区間 $I = (t_0, t_1)$ を inactive にする。そうでなければ, アクティブと呼ぶ。区間がアクティブならば, t_{01} を $t(t_0, t_1)$ の中央値(正確には, 中央値近似する S の要素)とする。 F_H 族が閉集合族なので, $P^{opt}(t_1) \supseteq P^{opt}(t_{01}) \supseteq P^{opt}(t_0)$ である。したがって, $P^{opt}(t_{01})$ の計算に際しては, H から $P^{opt}(t_0)$ の外側の頂点をすべて取り除き, $P^{opt}(t_1)$ の内部の頂点を単一頂点へ縮約できる。したがって, $|P^{opt}(t_0)| - |P^{opt}(t_1)|$ 個の頂点を持つ有向グラフ $H(I)$ の最大推移閉包を計算すれば十分である。

2 分分岐探索は各レベルにおいて, 前のレベルで作られたアクティブな区間を 2 つの区間へ分離する。初期の全区間の長さは高々 N であり, 葉区間の長さは少なくとも $1/(N-1) - 1/N = 1/N(N-1)$ であるので, レベルの数は $O(\log N^3) = O(\log N)$ である。各レベルで, すべてのアクティブな区間 I の $H(I)$ の頂点の数の合計は $O(n)$ である。大きさ n のグラフの最大推移閉包を求める計算時間は $T(n,n) = O(n^{1.5} \log n \log N)$ である。したがって, 各レベルを処理する計算時間は $O(n^{1.5} \log n \log N)$ である。したがって, 総計算時間は $O(n^{1.5} \log n \log^2 N)$ である。□

Theorem 6 与えられたセル c の直交領域 stabbed union の集合族に対する最適ピラミッドは $O(n^{1.5} \log n \log^2 N)$ 時間で計算できる。□

6. まとめと今後の課題

ピラミッド構造はいくつかの分野で有用である。まず, 確信度分布の傾向の良いビジュアル化を与える。次に, しきい値の高さ t 以上のピラミッドの一部の選択によって(このオペレーションを切り取りと呼ぶ), 領域内部のその実際幾何学的な位置に依存する各データに対するルールの影響に基づいた情報と一緒に結合ルールを与える領域を生成することができる。

一旦領域ルールを得れば, オリジナルのデータ分布からピラミッドを引き平均確信度をもって同じ先行のサポートを補充することにより, 容易にこのルールの影響を削除することができる。さらに, より弱いルールを抽出することができる。また, 異なるピークを備えたピラミッドを同時に考えることにより, データから 2 つ以上の階層化ルールを抽出できる。これにより,

各々のそのようなピラミッドからの高い一部分を切り取ることにより, データの大多数をカバーするクラスタリングを自動的に与える。第 4 に, すべてのピークでのピラミッド近似で離れ値であるデータアイテムを, 例外的なデータ(恐らくある入力エラー)あるいはあいまいなデータと見なすことができるので, データクリーニングのために潜在的に使用することができる。適切な領域族の選択は非常に重要な問題である。

データマイニングアプリケーションのための良い族を定義することに役に立ち, $d \geq 3$ (特に $d=3$) の時, 頂点の各出次数, 入次数の良いパラメーターを発見する実験を行う。

【文献】

- [AIM] R. Agrawal, T. Imielinski, A. Swami, Mining Association Rules between Sets of Items in Large Databases, *Proc. SIGMOD* (1993) 207--216.
 [CCKT02] D. Chen, J. Chun, N. Katoh, and T. Tokuyama, Layered Data Segmentation for Numeric Data Mining, Presented at Submitted.
 [CKT02] J. Chun, N. Katoh, and T. Tokuyama, How to Reform a Terrain into a Pyramid, Presented at *DIMACS Workshop on Geometric Graph Theory* (2002).
 [CST02] J. Chun, K. Sadakane, T. Tokuyama, Improved Algorithms for Constructing Pyramids from Terrains. Presented at *Japan Conference on Discrete and Computational Geometry* (2002).
 [FMMT96b] T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama, SONAR: System for Optimized Numeric Association Rules, *Proc. SIGMOD 1996* (1996) p.553.
 [FMMT99] T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama, Mining Optimized Association Rules for Numeric Attributes, *Journal of Computer and System Sciences* **58** (1999) 1-12.
 [FMMT01] T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama, Data Mining with Optimized Two-Dimensional Association Rules, *ACM Transaction of Database Systems* **26** (2001) 179-213.
 [YFMMT97] K. Yoda, T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama, Computing Optimized Rectilinear Regions for Association Rules, *Proc. KDD97* (1997) 96-103.

全 眞 嬌 Jinhee CHUN

東北大学大学院情報科学研究科システム情報科学専攻博士後期課程在学中。2003 年東北大学大学院情報科学研究科システム情報科学専攻博士前期課程修了。データマイニングに関する研究。電子情報通信学会学生会員, 日本データベース学会学生会員。

Danny Z. CHEN

Notre Dame 大学 Dep. of Computer Science and Engineering 教授。1992 年 Purdue 大学(Ph.D., Computer Science)。アルゴリズム, 計算幾何学, データマイニング等の研究に従事。ACM 会員, IEEE シニア会員。

加藤 直樹 Naoki KATOH

京都大学大学院工学研究科建築学専攻教授。1977 年 9 月京都大学大学院高学研究所博士後期課程中途退学(工学博士)。組合せ最適化, 離散アルゴリズム, 計算幾何学, 建築システム最適化, データマイニング等の研究に従事。情報処理学会, 電子情報通信学会, 建築学会, ACM 各会員。

徳山 豪 Takeshi TOKUYAMA

東北大学大学院情報科学研究科情報システム評価学分野教授。1985 年東京大学理学系大学院数学専門課程卒(理学博士)。理論計算機科学(特に計算幾何学を中心とした離散アルゴリズム理論), 離散数学, データマイニング等の研究を行う。情報処理学会アルゴリズム研究会主査, 電子情報通信学会, 日本数学会, 日本 OR 学会, 応用数理学会, ACM 各会員。