

Web ログ分析における対話的パターン検索のためのデータキューブモデル

6X-04

葛谷 雄一 大森 匡 星 守

電気通信大学大学院 情報システム学研究科*

1 研究の背景と目的

近年の DB 分野では, World Wide Web 等の情報サーバのアクセス記録データを大量に蓄積し, それらを分析して有用な情報を取り出す試みが盛んである. Web ログを分析する際によくあるパターン抜き出す手法として Apriori が用いられる [1]. しかし Apriori を用いて Web ログを分析すると多くの情報が出てきてしまい, 理解するのは困難であるため, 分析条件を変えて対話的に分析を繰り返す必要がある [2].

そこで, 本稿では Apriori を用いて Web ログを対話的に分析するための分析モデルとしてデータキューブに基づいた分析モデルを提案し, スライス, ロールアップに相当する演算をサポートすることで, 対話的に分析を行う方法とその演算の計算方法について述べる.

2 Web ログ分析

Apriori を用いて Web におけるよくあるアクセスパターンを抜き出すためには, Web ログにあるユーザーが一定時間の間にどのページを見たかというアクセスレコードの形に変換する必要がある. アクセスログからの変換は同一の IPaddress のログを時間を条件につなぎ合わせることで行なう.

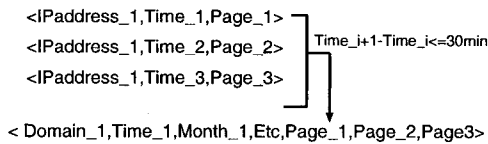


図 1: データ変換

本稿ではさらに Time を時間と月に変換し, IPaddress をドメインに変換し, その他 (Web 上でのイベントの有る無し) を追加し (図 1), 分析の対象として \langle ドメイン, 時間, 月, その他, 「ページ」 \rangle の 5 属性の形のアクセスレコードを用いる. (但し, ここで言う属性「ページ」は図 1 の様にページ集合を値にとる).

このアクセスレコードから良くあるパターンを抜き出すのに Apriori を用いる. しかし, 単純に Apriori を実行しても多くのパターンが出てきてしまい, どの情報が有用なのか分からない. そこで Apriori 実行時に制約条件をつけ, 何回も対話的に分析を繰り返す必要がある.

3 データキューブモデル

対話的に分析を繰り返すには, 出てくる情報の比較や, 制約条件の制御が簡単に行なえなければならない. そこで対話的に分析を繰り返すための分析モデルとしてデータキューブモデルを提案する (図 2).

3.1 用語説明

本稿で扱うアクセスレコードは属性の組みである. 属性は, 原子的な値を取る場合と, (「ページ」のように) 集合を値として取る場合がある. 一方, Apriori では, 命題変数をアイテムと呼び, アイテムから構成された (良く起きる) パターンを検出する. このパターン, すなわちアイテムの組合せをアイテムセットと呼び,良く起きるパターンを高頻度アイテムセットと呼ぶ. 任意の 1 アイテムはあるカテゴリに所属するとされる. 本稿では, レコードの属性をこの (アイテムの) カテゴリに対応させて考える. 即ち, Web ログの場合, ドメイン, 月, 「ページ」等が属性であり, あるドメイン d_i やあるページ p_j 等がアイテムになる (図 2a). 結果, アクセスレコードはそのレコードで成立するアイテムの列で表される.

アクセスレコードにおいて属性 X のアイテムが, 必ず 1 つしか存在しないならば, 属性 X をリレーション属性と呼び, 1 つ以上のアイテムが存在する可能性があるならば, 属性 X をアイテムセット属性と呼ぶ. Web ログの場合, ドメイン, 時間, 月, その他がリレーション属性であり, 「ページ」がアイテムセット属性である.

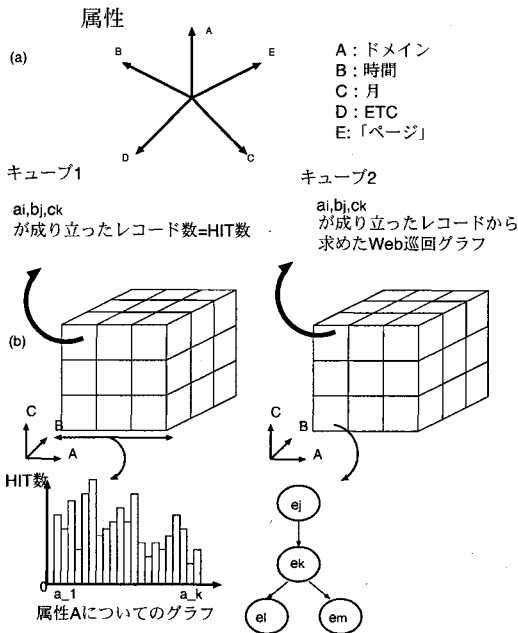
3.2 定義

ある属性の組みからなる空間を考えた時, 各属性を軸として構成される空間をキューブと呼ぶ. 図 2(b) は A, B, C を軸とした 3 次元のキューブである. キューブの樹目の一つ一つをセルと呼ぶ. セルは, 各属性のアイテム $a_i \in A, b_j \in B, c_k \in C$ が成立するレコード群から計算されたある値を持つ.

本稿では, キューブ 1 とキューブ 2 の 2 種類を用意する. キューブ 1 は, セルの値として, HIT 数を持つ. HIT 数とは, セルに対応するレコード群の総数のことである.

一方, キューブ 2 は, セルの値として, セル計算範囲のアクセスレコード群から求められた (アイテムセット属性「ページ」のアイテムからなる) 高頻度アイテムセットの集合を持つ. これから, Web 巡回グラフ (ユーザが Web をよくたどる経路) を計算できる. 以下では, キューブ 2 はセル値としてこのグラフを持つと見なし話を進める. また, キューブ 1 のセルの値はスカラーであるが, 軸

*Data Cube Model for Interactive Mining of Web-log analysis.
Y.TSUTATANI, T.OHMORI, M.HOSHI, Univ. Electro-Comm.



A, B, C についての3次元のキューブの場合

図 2: データキューブモデル

にアイテムセット属性が有る場合、従来のデータキューブとは異り単純に集計演算等を行なうことが出来ない。

3.3 演算

キューブ1, キューブ2 に対してスライスおよびロールアップ演算を適用することで効率的にユーザーの行動パターンの抽出と分類を行なうことができる。

・スライス

条件を指定して、指定された条件を満たす各セルの値を計算する。キューブ2をスライスする時はスライスされたセルの値を相互に比較するために、各セルの計算する情報の粒度を一定にする。(足切り値率 $\theta\%$ を揃える)。

・ロールアップ

ある属性のアイテムの幾つかを一つにまとめることでグループ化し、セルを新たに作り直す。

3.4 分析手順の例

キューブ1の値を見ること(図3(a))で大まかな分析の方針を立て、ヒット件数の多いドメインと少ないドメインにグループ化(図3(b))して、それぞれのキューブ2の値の計算を行なってグループ毎の行動の違いを見たり、さらに月でスライスすることでそれらのグループの行動を時間軸に沿って見ることが出来る(図3(c))。このようにスライスとロールアップを繰り返すことで図3の様な結果を得ることが可能である。

3.3節で説明したスライスは、図3(a)や(c)を各々計算することである。ロールアップは図3(b)の様にグルー

(a) Aのみから成るキューブ1

HIT	127	156	17	327	278	212	32
A	a1	a2	a3	a4	a5	a6	a7

(b) A上のロールアップ

- a3, a7 → G1
- a1, a2 → G2
- a4, a5 → G3

(c) A, Cから成るキューブ2

G1	←	←	↖
G2	↓	↗	↗
G3	↖	↗	↗
	4月	5月	6月

図 3: 分析結果一例

プG1, G2, G3 を作成し分析の階層を上げることである。

3.5 アルゴリズム

リレーション属性のみでスライスする場合、セルの値を計算する時の各セルに対応したレコードの一つ一つは重複しないが、アイテムセット属性でスライスする場合は重複する可能性があるため、それぞれ個別にセルの値を計算するとデータベースのスキャンに無駄が生じ、候補アイテムの数上げにも重複が生じる。そこでアイテムセット属性のアイテム p_i ($1 \leq i \leq N$) に対応したセル i のキューブ2の値を足切り値率 $\theta\%$ で計算する際のアルゴリズムを説明する。

[手順1 候補アイテムセット生成]

あるアイテムセット I が候補アイテムセットであるためには I の部分集合の全てがあるセル i で高頻度アイテムセットでなければならない

[手順2 候補アイテムセットの各セル毎のサポート数の数え上げ]

p_i ($1 \leq i \leq N$) の内一つでも成り立つアクセスレコード群に対して、あるレコードで p_i が成り立つならば、セル i について候補アイテムセットのサポート数の数え上げる

[手順3 高頻度アイテムセットの生成]

あるセル i で候補アイテムセット I が高頻度であるためには I のサポート数が $hit_i \times \theta/100$ 以上でなければならない (hit_i はキューブ1のセル i の値)。手順1に戻る。

4 まとめ

本稿では、対話的な分析を行なう為の問題点を整理し、解決の方法としてデータキューブモデルを提案し、さらにスライスの高速化のアルゴリズムを提案した。

今後の課題として得られた情報の統計的な手法への適用と、ロールアップの計算の高速化等を考えている。

参考文献

- [1] R. Agrawal, et al. Fast Algorithms for Mining Association Rules, Proc.20th VLDB, pp.487-499, 1994.
- [2] U. Dayal, Web Content and Usage Mining for E-Commerce Applications, Proc. DBWEB2000, pp.311-316, 2000.