

データマイニングにおける追加データの処理方式

2P-4

大森 匡 Zhang Ji 星 守  
電気通信大学 大学院 情報システム学研究科\*

1 はじめに

最近、データベースを用いた連想型ルールのマイニングが注目されており、Apriori（文献1）などの高速な算法が提案されている。しかし、これらの算法はいずれも静的な（変更されない）データベースを対象としている。一方で、実際のデータベースではデータが増加していくことが普通である。本稿では、こうした「増加していくデータベース」において Apriori 法を効率良く適用する方式を提案する。

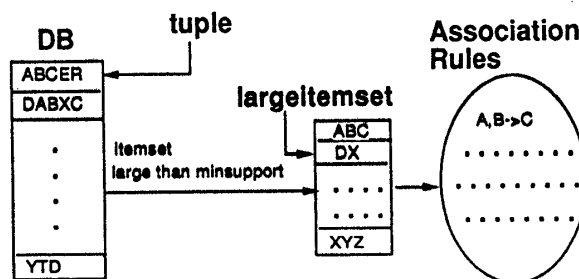


図 1: マイニングの流れ

2 ルール発見問題

はじめに連想型ルール発見問題と Apriori 法を図 1 を用いて説明する。

今、データベース DB のスキーマを  $[A_1, A_2, \dots, A_n]$  として各属性がある事象の発生の有無を示すブール値としよう。すなわち、DB のタプルは「命題  $A_1, A_2, \dots, A_n$  のうち どれが同時に成立したか」を示すログレコードであり、DB は大規模なログデータベースである。（図 1 ではタプルは「同時に成立した属性のリスト」で表現されている。例えば、タプル ABCER は「事象 A, B, C, E, R が同時に発生した」という記録を表す。）また、以下では DB における属性組  $Ax Ay \dots Az$  の出現回数（サポート数）とは、「DB において属性  $Ax, Ay, \dots, Az$  が全て 1 であるようなタプルの総数」のことである。このとき、連想ルール発見問題とは、サポート数が  $\theta$  以上となるような属性組 (ItemSet) を DB に対し全て求め、そこから「A, B なら C である確率は  $\alpha$ 」という形のルールを発見することである（図 1）。（ $\theta$  は minsupport と呼ばれ、minsupport より出現回数の大きい属性組を Large ItemSet、それ以外の属性組は Small ItemSet と呼ぶ。ルール自体は Large ItemSet から条件つき確率の式で求められる。）

一般にルール発見問題では  $k$  個の属性からなる属性組についてそのサポート数を DB をスキャンして数え、 $k = 1$  から順に何回か DB をスキャンして全ての Large Itemset を見つける。Apriori 法では、「ある属性組  $ItemSetX$  の中で、もし、一つかそれ以上の部分集合が Small Itemset であれば、 $ItemSetX$  自体も Small Itemset である」という原則を使い、数えるべき属性組を限定していく。

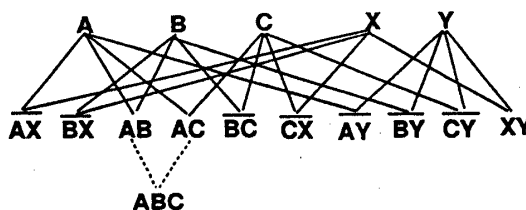


図 2: Apriori

図 2 に Apriori が終了した時の属性組の例を示す。図中、上に線がついている属性組は Small ItemSet であり、線がついていない属性組は Large ItemSet である。属性組の間の直線は、上位の属性組 ItemSetX と ItemSetY を結合 (i.e. 和集合をとること) して得られた属性組 ItemSetXY について、それを DB スキャンの時に数えることを示している。点線は「上位の属性組を結合して得た属性組  $IS1$  の部分集合の中に Small ItemSet が入っているので  $IS1$  は数えない」を示す。（図 2 では、属性組 ABC は BC が small だから数えない。また、AX, BY などは small であるため、これ以上結合しない。）このように Apriori ではサポート数を数えた Small Itemset に対し、それを真に含む属性組は数えられない。

3 追加データの処理方式

それではデータベース DB に対し、タプルを追加する ( $DB_{new} := DB_{old} + \Delta DB$ ) した場合を考えよう。このとき、 $DB_{old}$  においては Small Itemset だった属性組が  $DB_{new}$  では Large Itemset になりうるため、これを見つけないといけない。特に、Apriori では一旦 Small ItemSet と判断するとそこから後の属性組のサポート数

\*Mining Association Rules from Incremental Database  
J.Zhang, T.Ohmori, M.Hoshi (U.Electro-Comm.)

を数えないわけで、追加データ処理のときはこうした DBold における Small ItemSet  $X_1X_2\dots X_n$  を DBnew において数え直す必要がある。このとき、DBold を全スキャンするのは増分が非常に微量であれば不利である。すなわち、元のデータベースを前処理しておき、新しい Large Itemset を見つける時にはスキャンするデータベース領域を限定すべきである。以下では数え直しのアクセスボタンに応じて DBold を分解してアクセス量を減らす方法を提案する。これは DBold への前処理方法 (3.1 節の圧縮と DB 分割) とデータ追加時の処理 (3.2 節) とにわけられる。

### 3.1 前処理 - DB の圧縮と分割

#### 前処理 1 - DB の圧縮と出現回数の記憶

最初に Apriori を DB に適用した際、次の情報を記憶しておく。

1: 次のような {Large Itemset  $LX$ , Counter} の組を記録しておき、DB から余分なタプルを削除する。すなわち、Counter とは DBold において Large Itemset  $LX$  に含まれる属性にのみ 1 が立っているようなタプルの数である。この Counter が記録されると、これらのタプルは DBold から除いても問題ない。DBold において残ったタプル集合を RDB (Remaining DB) としておく。

2: この他に、DBold においてサポート数を数えた結果 Small Itemset と判断された属性組についても [属性組, サポート数] を記憶しておく。

#### 前処理 2 - Signature による DB 分割

次に RDB をいくつかのグループに分割し、追加データ処理のときに RDB のアクセス領域を限定する。一般に追加データ  $\Delta DB$  が生じた時、 $\Delta DB$  を一度スキャンして前処理 1 の情報を使うと新しい Large ItemSet が見つかる場合がある。問題となるのは、この新しい Large Itemset  $NewLx$  が DBold では Small ItemSet であったため、 $NewLx$  を真に含む属性組  $X_1\dots X_n$  を DBold で数えていないことである。実際、 $X_1\dots X_n$  のサポート数は DBold で数え直さないといけない。そこで、RDB へのアクセス量を削減するために RDB をシグニチャを使って分割しておく。すなわち、

1: タプル  $t$  で 1 が立つ属性を  $Y_1\dots Y_n$  としたとき、そのシグニチャを  $\text{sig}(t) = h(Y_1) + \dots + h(Y_n)$  とする。DB の属性数を  $N$  とすると、 $t$  は  $N$ -bit 長レコードであり、signature は  $k$ -bit 長で生成する。(  $h$  は  $k$ -bit 長レコードを生成するハッシュ関数で、 $+$  はビット単位の OR 演算。)

2: 次に、RDB を  $\text{sig}(t)$  の値に応じて  $2^k$  個のクラスタへ分割しておく。(各クラスタは同じ  $\text{sig}(t)$  値を持つタプル集合。)

こうすると、追加データ処理時に属性組  $X_1\dots X_n$  のサポート数を RDB で数えるときには  $\text{QueryMask}(X_1\dots X_n) = h(X_1) + \dots + h(X_n)$  を計算し、QueryMask にヒットするシグニチャを持つクラスタ (QueryMask で 1 が立つビット位置に 1 が立っているようなクラスタ) だけをアクセスすれば良い。

この利点は、RDB を互いに素なクラスタへ分割しておけることである。特に属性間で相関度が強い場合、シグニチャの生成方法を Large ItemSet の出方に応じて決めるとうまくいきそうである。

極めて単純な例をあげると、図 2 において DB の全属性が  $[A, B, C, X, Y]$  に対し極大な Large ItemSet は  $AB$ 、 $AC$ 、 $XY$  であるから  $\text{sig}(t)$  を  $[t$  が  $AB$  を含むか、 $AC$  は?,  $XY$  は?] の 3-bit 長シグニチャとして、RDB を [000] から [111] までの 8 クラスタへ分解しておくことが考えられる。こうすると、属性組  $ABC$  の数え直しに際してはクラスタ [110] だけをアクセスすれば良い。

一般にデータベースにおいては属性の総数  $N$  が  $10^5$  程度と非常に大きい、各タプルの 1 が立つビット数は 10 程度と小さいためシグニチャ化するには都合が良い。本節で提案した方式ではシグニチャのビット数  $k$  として  $2^k$  個に DB をクラスタ分割するが、 $k = 10$  としてクラスタ総数が 1000 個でヒット率 1% 程度なら見合うというものである。

### 3.2 追加データに対する処理

以上に述べてきた前処理の後、追加データ  $\Delta DB$  に対し次の手順で処理を行なう:

手順 1:  $\Delta DB$  をスキャンし、前の (DBold における) Large Itemset に対しサポート数の増分を計算し DBnew でのサポート数を得る。

手順 2: DBold においてサポート数を数えた Small Itemset に対し、 $\Delta DB$  をスキャンしてそのサポート数を更新する。これにより新しい Large Itemset が見つかる。

手順 3: 手順 2 で見つかった新しい Large ItemSet を各々  $NewLx_1, \dots, NewLx_N$  とする。これら属性組から始めて  $RDB + \Delta DB$  に対し Apriori を実行する。このとき、ある属性組  $X_1\dots X_n$  を DBold で数えていなかった場合、前処理 2 で作成した RDB のクラスタに対し  $\text{QueryMask}(X_1\dots X_n)$  を作成してこれにヒットしたクラスタを使って Apriori を実行する。

手順 4: 追加データに対する処理を完了した後、次回処理のために  $DBold + \Delta DB$  に前処理 1、2 を行なう。

## 4 まとめ

本稿では連想型ルール発見において追加データを扱う方法を提案した。具体的には、シグニチャを使って元のデータベースを分割することで Apriori の再計算負荷を削減した。シグニチャの作り方は Large ItemSet の出方、つまりデータベースの分布に応じた設計項目であり、現在試験的な評価を行なっている。

参考文献 1: Agrawal 他. "Fast Algorithms for Mining Association Rules". Proc. 20th VLDB 94.