

Hellinger 엔트로피를 이용한 다차원 연속패턴의 생성방법

이 창 환^{*}**요 약**

데이터 마이닝에서 연속패턴(sequential pattern) 생성기술은 시차를 두고 발생한 사건들에 대하여 잠재해있는 패턴을 발견하는 기술을 의미 한다. 본 연구는 정보이론을 이용하여 데이터베이스로부터 연속패턴을 자동으로 발견하는 방법에 관한 내용이다. 기존의 방법들이 한 속성내에서의 연속패턴만을 탐지하는 일차원 연속패턴을 생성하는데 비하여 본 연구에서 제시하는 방법은 데이터베이스내의 모든 속성간의 연속패턴 관계를 탐지할 수 있는 다차원 연속패턴을 생성할 수 있다. 본 연구에서는 연속패턴 생성을 위하여 헬링거(Hellinger) 변량을 사용하였으며 이를 이용하여 발견된 연속패턴들의 중요도를 측정할 수 있었다. 또한 헬링거 변량의 함수적인 특성을 분석하여 연속패턴 추출의 복잡도를 줄이기 위한 두 가지의 법칙이 제안되었고 다수의 실험 데이터를 통하여 다차원의 연속패턴을 생성할 수 있음을 보였다.

Learning Multidimensional Sequential Patterns Using Hellinger Entropy Function

Chang-Hwan Lee^{*}**ABSTRACT**

The technique of sequential pattern mining means generating a set of inter-transaction patterns residing in time-dependent data. This paper proposes a new method for generating sequential patterns with the use of Hellinger measure. While the current methods are generating single dimensional sequential patterns within a single attribute, the proposed method is able to detect multi-dimensional patterns among different attributes. A number of heuristics, based on the characteristics of Hellinger measure, are proposed to reduce the computational complexity of the sequential pattern systems. Some experimental results are presented.

키워드 : 데이터마이닝(Data Mining), 연속패턴(Sequential Pattern), 기계학습(Machine Learning)

1. 서 론

인간의 직관으로부터 직접적으로 얻어진 규칙은 정보 전달의 어려움 등으로 인하여 불합리하고 많은 비용이 듈다. 그렇기 때문에 주어진 데이터베이스로부터 전문가 수준의 자동적으로 추론되어진 규칙을 사용할 수 있다면 대단히 유용 할 것이다[1]. 데이터 마이닝은 대량의 데이터로부터 의미있는 패턴을 자동으로 발견하는 기술로서 현재 많은 연구가 진행되고 있다[2]. 본 연구는 데이터 마이닝의 기술 중에서 연속패턴의 생성에 대한 연구로서 시계열 데이터로부터 잠재해 있는 연속패턴을 정보이론을 이용하여 IF-THEN의 형식으로 자동 생성하는 기술에 관한 연구이다.

연속패턴의 생성에 대한 연구는 현재 몇 가지의 방법이 제

안되어 있으며 상용화되어서 사용하고 있는 알고리즘도 있다. 이를 중에서 현재 가장 널리 사용되고 있는 연속패턴의 방법은 IBM에서 개발한 GSP 알고리즘[3]이다. GSP 알고리즘은 역시 IBM에서 개발하여 널리 사용하고 있는 연관(association) 알고리즘[3-5]의 기능을 확장하여 연관패턴을 생성할 수 있게 하였다.

GSP 알고리즘의 연속패턴 생성기술은 고객이 제품을 구입한 데이터를 분석하여 각 제품의 구매형태가 시차에 따라서 어떠한 패턴을 보이는지를 발견하는 기술로서 이 기술들은 구체적으로 고객의 구매 기록 데이터를 분석하여서 다음과 유사한 법칙들을 생성한다.

IF 제품 = TV, THEN 제품 = VCR

이 법칙의 의미는 TV를 구매한 고객은 추후에 VCR을 구

* 종신회원 : 동국대학교 정보통신학과 교수
논문접수 : 2003년 11월 21일, 심사완료 : 2004년 4월 6일

입한다는 의미이다. 이러한 기준의 연속패턴 생성의 방법들은 IF 부분에 오직 한 개의 속성값을 허용하며 THEN 부분에는 IF 부분에서 사용된 속성만이 사용될 수 있다. 또한 이 계열의 많은 알고리즘들은 시간창(time window)을 두어서 주어진 시간간격 이내에 발생한 데이터만을 고려하여 법칙을 생성한다.

하지만 이러한 기준의 연속패턴 방법들은 전체 데이터 중에서 특정한 한 개의 속성(대부분 제품)에 대해서만 시간에 따라 발생하는 연속패턴을 탐색·생성하는 법칙은 훨씬 의미있고 광범위한 정보를 제공하며 훨씬 많은 응용 범위를 가진다.

이와 같은 이유로 기준의 연속패턴 방법은 대량의 제품을 취급하는 업체에서는 일부 분석효과를 기대할 수 있지만 판매 제품이 소수인 경우에는 기준의 연속패턴 방법의 적용은 그 의미를 찾을 수 없다. 즉 판매 제품의 종류가 작을 경우에 이들 제품간의 판매 연관성보다는 이들 제품을 각각 구입하게 하는 다른 속성의 원인 분석이 더욱 중요한 정보가 될 것이며 따라서 한 개의 속성 내에서의 연속패턴을 제공하는 기능은 소수 제품의 환경에서는 거의 의미를 찾을 수 없다.

본 연구에서 제시하는 연속패턴 생성 방법은 테이블내의 모든 속성들의 값에 대하여 서로의 연속패턴 연관관계를 계산할 수 있다. 따라서 훨씬 중요한 의사결정의 정보를 제공하며 훨씬 다양한 분야에 대하여 적용할 수 있다. 예를 들어서 본 연구의 법칙생성의 방법을 수행하면 다음과 같은 연속패턴의 법칙을 생성할 수 있다.

- 1) IF 제품=SOAP, THEN 제품=SHAMPOO
- 2) IF 직업=회사원 & 제품=TV, THEN 제품=VCR

위에서 법칙 1)의 의미는 SOAP를 구입한 사람은 다음에 SHAMPOO를 구입한다는 의미로서 이는 기준의 연속패턴 생성알고리즘의 기능과 동일한 기능을 보여준다. 법칙 2)의 의미는 직업이 회사원인 사람이 TV를 구입하면 다음에 VCR을 구입한다는 의미이다.

따라서 본 논문의 방법에서 생성하는 법칙은 훨씬 의미 있고 광범위한 정보를 제공하며 훨씬 많은 응용 범위를 가진다.

2. 본 연구의 내용

본 연구에서 제시하는 방법은 데이터를 읽어서 다음과 같은 형식의 연속패턴 법칙을 생성하는 내용이다. 예를 들어

$B_{i_1}, B_{i_2}, B_{i_3}$ 을 테이블의 속성이라고 하고 $b_{j_1}, b_{j_2}, b_{j_3}$ 를 각 속성에 대한 임의의 값으로 가정하자. 이때 본 알고리즘에서 생성되는 법칙은 다음과 같은 형식으로 표현된다.

$$\text{IF } B_{i_1} = b_{j_1} \wedge B_{i_2} = b_{j_2} \wedge B_{i_3} = b_{j_3} \cdots, \text{THEN } A = a_k$$

이 법칙의 의미는 IF 부분에 있는 $B_{i_1} = b_{j_1} \wedge B_{i_2} = b_{j_2} \wedge B_{i_3} = b_{j_3} \cdots$ 의 조건을 만족하는 행위가 발생하였을 때 그 후에 $A = a_k$ 를 만족하는 행위가 발생하는 패턴이 있다는 의미이다.

이와 같은 연속패턴 법칙 생성의 기본적인 가정은 IF 부분의 조건이 THEN 부분의 확률 분포에 영향을 준다는 가정에서 출발한다. 직관적으로 말하면, 특정 속성(attribute)의 값이 정해질 때 목표 속성(target attribute)의 확률 분포를 현저히 변화시킨다면 이는 특정 속성의 값을 결정하는 중요한 역할을 의미한다. 따라서 시스템은 우선 속성의 변수값(예를 들어서 속성 B 가 b 의 값을 가짐)을 우선 선택하고, $B = b$ 이라는 사건이 목표 속성 A 의 분포값에 어떤 영향을 끼치는지를 점검한다. 만약 그것이 목표 속성의 확률 분포에 상당한 영향을 끼친다면 시스템은 다음과 같은 규칙이 있음을 가정한다[7].

$$\text{IF } B = b, \text{THEN } A = a \text{ with } H$$

각 규칙은 각 규칙의 중요도를 표현하는 H 값을 포함하고 있다. 여기서 THEN 부분에는 오직 한 개의 속성값이 나타나도록 하였고 규칙의 왼쪽 IF 부분은 논리곱(conjunction)의 형태로 되어 있다. 또한 논리합(logical OR)과 부정 논리(logical NOT)의 형태는 고려하지 않는다.

본 연구에서는 이와 같은 목표 속성의 확률분포의 변화 정도를 측정하기 위하여 헬링거(Hellinger) 엔트로피 함수[8]를 사용하였다. 본 연구의 연속패턴 환경에서 위와 같은 법칙의 경우 헬링거 함수는 다음과 같이 정의된다.

$$[\sum_{a \in A} (\sqrt{P(a)} - \sqrt{(P(a|b))^2})^2]^{\frac{1}{2}}. \quad (1)$$

위의 정의는 속성 A 에 대해서 $B = b$ 의 값 이전의 확률 분포와 $B = b$ 값 이후의 확률분포에 대하여 얼마나 차이가 있는가에 대한 계산식이다.

연속패턴 생성의 법칙에서 고려해야 할 사항은 목표 속성의 값 중에서 가장 빈번히 나타나는 하나의 값은 규칙의 오른 편에서 나타나고 다른 값들은 확률 $1 - P(a)$ 에 포함되어지는 것이다. 다시 말해 규칙의 정확도를 측정하는 데에 있

어서 중요한 점은 특정 데이터가 규칙의 오른 편의 목표값 ($A = a$)에 속하는지를 검사하는 것이다. 예를 들어서 목표 속성 A 가 k 개 (a_1, a_2, \dots, a_k)의 값을 가지고 있고 연속패턴 규칙은 오른 편에서 $A = a_1$ 이라고 가정하자. 그러면 속성 A 의 확률 분포가 $(P(A = a_1), P(A \neq a_1))$ 와 같이 2진의 확률 분포로 변환된다. 그래서 식 (1)은 다음과 같이 변환될 수 있다.

$$(\sqrt{P(a|b)} - \sqrt{P(a)})^2 + (\sqrt{1-P(a|b)} - \sqrt{1-P(a)})^2 \quad (2)$$

이 식에서 $P(a|b)$ 는 $B = b$ 라는 조건 하에서 $A = a$ 의 조건 확률을 의미한다.

구체적으로 얘기하면 먼저 IF-THEN 법칙에서 THEN 부분에 나타나는 목표 속성에 대한 이전확률분포 분포를 구한다. 그 다음으로 IF 부분의 조건을 만족한 상태에서의 목표속성에 대한 확률분포인 이후확률분포를 계산한다. 이후 확률분포를 계산할 때는 IF 부분의 행위가 THEN 부분의 행위보다 먼저 시행된 데이터의 개수만 고려를 하여 계산한다. 이와 같은 방식으로 이전확률분포 와 이후확률분포를 계산한 후에 이들이 서로 얼마나 상이한가의 정도를 헬링거 엔트로피 함수를 사용하여 계산한다. 이와 같이 계산된 엔트로피 함수의 값이 해당 법칙의 정확도를 의미한다.

다음과 같은 연속패턴 법칙에 대한 정확도를 계산한다고 가정하자.

IF $A = a \wedge B = b$, THEN $C = c$

이 경우 속성 C 의 이전확률분포를 계산하는 것은 법칙 생성 알고리즘의 경우와 동일하다. 하지만 이후확률분포를 계산할 때 본 알고리즘은 전체 데이터 중에서 조건 $A = a \wedge B = b$ 를 만족하는 각 데이터마다 해당 데이터의 행위자(보통 고객임) 데이터가 끝날 때까지의 잔여 레코드 중에서 조건 $A = a \wedge B = b$ 를 만족하는 모든 레코드를 원래의 데이터에 추가되는 것으로 가정하고 이렇게 수정된 데이터의 분포를 이용하여 이후확률분포를 계산한다.

구체적인 예를 들어서 (그림 1)의 경우 다음의 법칙에 대한 정확도를 계산해 보자.

IF 고객 = C1 \wedge 제품 = P1, THEN 제품 = P2

이 법칙의 목표속성은 제품이며 제품에 대한 이전확률분포는 쉽게 계산할 수 있다. 그 다음으로 이후확률분포를 계산할 때 먼저 IF 부분의 조건 고객 = C1 \wedge 제품 = P1을 만족하는 레코드는 (1), (3), (7)임을 알 수 있다. 따라서 레코드 (1)에 의하여 레코드 (3)과 (7)이 원래의 데이터에 추가

되는 것으로 가정하며 레코드 (3)에 의하여 레코드 (7)이 추가되는 것으로 가정한다. 즉 레코드 (7)은 이 법칙에 대하여 3번 추가된다. 이와 같이 수정된 데이터의 분포를 이용하여 이후확률분포를 계산한다. 그리고 이러한 이전확률분포와 이후확률분포의 차이를 헬링거 합수로 측정한 값이 위의 법칙에 대한 정확도가 된다.

고객	일자	제품	
C1	D1	P1	(1)
		P4	(2)
	D2	P1	(3)
		P2	(4)
		P3	(5)
	D3	P2	(6)
		P1	(7)
...

(그림 1) 연속패턴 생성의 예제

Input : IF $B_1 = b_1 \wedge B_2 = b_2 \wedge \dots \wedge B_k = b_k$ THEN $A = a$
Output : posterior prob. distribution of attribute A ;

```

read entire data file  $D$  and calculate prior distribution of attribute  $A$  ;
for each record  $R$  in  $D$  do
    if  $R$  satisfies  $B_1 = b_1 \wedge B_2 = b_2 \wedge \dots \wedge B_k = b_k$ 
        then
            PERSON := the value of attribute person of the record
            do
                if the current record satisfies
                     $B_1 = b_1 \wedge B_2 = b_2 \wedge \dots \wedge B_k = b_k$ 
                    then calculate the posterior dist. of attribute  $A$ 
                        using the method described in page 5 ;
                        read next record in  $D$  ;
                while (PERSON = the value of attribute person)
                    else
                        continue ;
                    end-if
                end-for
            return (the posterior prob. distribution of attribute  $A$ )

```

(그림 2) 연속생성에서의 이후확률분포 계산코드

연속패턴 생성에서 우리가 고려해야 할 또 다른 문제는 규칙이 얼마나 일반적인(general) 규칙인가 하는 것을 결정하는 것이다. 규칙을 만족하는 데이터의 수가 많을수록 그 규칙은 더욱 일반성이 있다고 할 수 있다. 그래서 규칙의 중요도의 부분으로서 규칙의 보편성의 정도를 계산하여야 한다. 규칙의 일반성을 나타내기 위하여 본 시스템에서는 다음 수식을 사용하였다.

$$\sqrt{P(a)} \quad (3)$$

결과적으로 H 계산은 규칙의 중요도를 주어진 규칙의 정확도와 일반성의 복합적인 계산으로 표현한다. 주어진 규칙에서 중요도의 마지막 형태는 다음과 같이 규칙의 일반성과 정확도 사이에서 생성된 식으로 정의된다.

[정의 1] 다음과 같은 규칙

$$\bigwedge_i B_i = b_{ij} \rightarrow A = a_k$$

의 H 값은 다음과 같이 정의된다.

$$H = \sqrt{P(\bigwedge_i B_i = b_{ij})} \\ \left[\left(\sqrt{P(A = a_k | \bigwedge_i B_i = b_{ij})} - \sqrt{P(A = a_k)} \right)^2 \right. \\ \left. + \left(\sqrt{(1 - P(A = a_k | \bigwedge_i B_i = b_{ij})} - \sqrt{(1 - P(A = a_k))} \right)^2 \right]$$

위의 식에서 B_i 는 i 번째 속성을 나타내며 b_{ij} 는 속성 B_i 의 j 번째 값이다.

3. 연속패턴생성의 방법

연속패턴생성 알고리즘을 수행하기 위해서 데이터는 (그림 1)과 같이 행위자(제품 판매의 경우에는 고객)와 시간 및 목적물(제품 판매의 경우에는 제품)의 속성을 포함하고 있어야 한다. 또한 알고리즘의 수행을 위해서 데이터는 행위자와 시간의 순서대로 정렬(sorting)을 시켜야 하며 연속형 속성은 이산 속성(discrete attribute)의 형태로 변환되어야 한다. 본 논문의 연속패턴생성 알고리즘의 전체적인 기능은 위와 같이 정리된 데이터의 값을 읽고서 k 개의 가장 의미있는 연속패턴법칙을 생성하는 알고리즘이다.

본 시스템에서 규칙을 생성하는 방법을 간략히 설명하면 아래와 같다. 먼저, 규칙의 원편이 한 개의 속성조건만을 갖는 단일 조건(single-condition) 규칙들을 생성한다. 알고리즘은 이를 단일 조건 규칙들에서 출발하여 가능한 왼쪽 면을 통한 깊이우선(depth-first) 탐색을 수행한다. 단일 조건들 중에서 H 계산값이 가장 높은 k 개의 규칙들이 규칙리스트의 형태로 저장된다. 이들의 H 값 중 가장 작은 H 값은 H_s 로 정의된다. 이 규칙 리스트의 각 원소에 대하여 알고리즘은 더욱 세분화된(specialized) 규칙들을 생성하려고 한다. 즉 알고리즘은 규칙 리스트에서 한 원소를 뽑고 추가적인 속성 조건을 원편에 추가하여서 세분화시킨다. 알고리즘은 다음 절에서 설명하는 정리를 중에서 하나를 만족할 때까지 세분화를 계속한다. 이 과정을 더 이상 세분화시

킬 법칙이 없을 때까지 반복하고 최종적으로 k 개의 가장 H 값이 높은 법칙을 출력한다.

3.1 H 값을 이용한 가지치기 방법

연속패턴 생성 시스템이 다른 데이터는 속성들의 숫자가 증가함에 따라 계산해야 할 규칙들의 총 숫자는 폭발적으로 증가한다. 예를 들어서 데이터가 r 개의 속성을 가지고 각 속성들은 v_i ($i = 1, \dots, r$)의 값들을 가진다고 가정하자. 그러면 최악의 경우 총 규칙의 개수는 다음과 같이 주어진다.

$$\prod_{i=1}^{r-1} (v_i + 1) - 1$$

우리는 헬링거 함수의 특성을 이용한 가지치기 기술을 제시한다. 우선 다음과 같은 규칙을 가정하자.

$$\text{IF } B = b, \text{ THEN } A = a \quad (4)$$

또한 $C = c$ 라는 조건을 더해서 다음과 같이 특수화된 규칙을 가정하자.

$$\text{IF } B = b \wedge C = c, \text{ THEN } A = a \quad (5)$$

식 (6)과 식 (7)에서의 규칙들을 각각 Rg 와 Rs 라고 하자. 규칙 Rg 와 Rs 의 H 값을 각각 H_g 와 H_s 라고 가정하면 이들은 다음과 같이 정의된다.

$$H_g = \sqrt{P(b)} [(\sqrt{P(a|b)} - \sqrt{P(a)})^2 \\ + (\sqrt{1 - P(a|b)} - \sqrt{1 - P(a)})^2] \\ = \sqrt{P(b)} [2 - 2\sqrt{P(a|b)P(a)} \\ - 2\sqrt{(1 - P(a|b))(1 - P(a))}] \quad (6)$$

$$H_s = \sqrt{P(bc)} [2 - 2\sqrt{P(a|bc)P(a)} \\ - 2\sqrt{(1 - P(a|bc))(1 - P(a))}] \\ = \sqrt{P(c|b)} \sqrt{P(b)} [2 - 2\sqrt{P(a|bc)P(a)} \\ - 2\sqrt{(1 - P(a|bc))(1 - P(a))}] \quad (7)$$

이 경우, 속성 C 의 값에 관계없이 우리는 다음과 같은 결과를 얻을 수 있다.

[정리 1] 목표 속성의 클래스 개수를 m 이라고 할 때 H_s 값은 다음의 경계값을 초과할 수 없다.

$$H_s \leq \max \{ \sqrt{P(a|b)} \sqrt{P(b)} [2\sqrt{m} - 2\sqrt{P(a)}], \\ 2\sqrt{P(a)} - \sqrt{(1 - P(a|b))\sqrt{P(b)}} \\ [2\sqrt{P(a)} + 2\sqrt{(1 - P(a))}] \}$$

[증명] 위의 주어진 법칙을 Rg 세분화된 법칙을 Rs 라고 하고, 이들의 H 값을 각각 Hg 및 Rs 라고 하자. 그러면 Hg 와 Rs 의 값은 식 (6)과 식 (7)과 같이 주어진다.

그러면 다음 수식이 성립한다.

$$\begin{aligned} P(ab) &= P(abc) + P(ab \neg c) \\ &= P(a|bc)P(bc) + P(a|b \neg c)P(b \neg c) \\ P(a|b) &= P(a|bc)P(c|b) + P(a|b \neg c)P(\neg c|b) \\ &= P(a|bc)P(c|b) + P(a|b \neg c)(1 - P(\neg c|b)) \\ \text{따라서 } P(c|b) &= \frac{P(a|b) - P(ab \neg c)}{P(a|bc) - P(ab \neg c)} \end{aligned} \quad (8)$$

위의 식에서 $\omega = P(ab \neg c)$ 라고 하자.

i) 첫째 경우, ω 의 값이 $P(a|b)$ 및 $P(a|bc)$ 보다 작은 경우
($\omega \leq P(a|b)$, $\omega < P(a|bc)$), $\max_{\omega} P(c|b) = \frac{P(a|b)}{P(a|bc)}$
이 되며 이때 ω 의 값은 0이 된다.

따라서

$$\begin{aligned} H_s &\leq \sqrt{\frac{P(a|b)}{P(a|bc)}} \sqrt{P(b)} [2 - 2\sqrt{P(a|bc)P(a)} \\ &\quad - 2\sqrt{(1 - P(a|bc))(1 - P(a))}] \\ &= \sqrt{P(a|b)} \sqrt{P(b)} [\frac{2}{\sqrt{P(a|bc)}} - 2\sqrt{P(a)} \\ &\quad - 2\sqrt{\frac{1 - P(a|bc)}{P(a|bc)}(1 - P(a))}] \\ &\leq \sqrt{P(a|b)} \sqrt{P(b)} [\frac{2}{\sqrt{P(a|bc)}} - 2\sqrt{P(a)}] \end{aligned}$$

본 알고리즘은 목표속성의 값 중에서 가장 빈도수가 높은 값을 만을 법칙의 THEN 부분에 포함시키므로 우리는 $P(a|bc)$ 값의 범위를 $1/m \leq P(a|bc) \leq 1$ 로 가정할 수 있다.

따라서 $H_s \leq \sqrt{P(a|b)} \sqrt{P(b)} [2\sqrt{m} - 2\sqrt{P(a)}]$

ii) 둘째의 경우, ω 의 값이 $P(a|b)$ 및 $P(a|bc)$ 보다 큰 경우
($\omega > P(a|b)$, $\omega > P(a|bc)$), $\max_{\omega} P(c|b) = \frac{1 - P(a|b)}{1 - P(a|bc)}$
이 되며 이때 ω 의 값은 1이 된다.

따라서

$$\begin{aligned} H_s &\leq \sqrt{\frac{1 - P(a|b)}{1 - P(a|bc)}} \sqrt{P(b)} [2 - 2\sqrt{P(a|bc)P(a)} \\ &\quad - 2\sqrt{(1 - P(a|bc))(1 - P(a))}] \end{aligned}$$

$$\begin{aligned} &= 2\sqrt{\frac{1 - P(a|b)}{1 - P(a|bc)}} \sqrt{P(b)} - \sqrt{1 - P(a|b)} \sqrt{P(b)} \\ &\quad [2\sqrt{\frac{P(a|bc)}{1 - P(a|bc)}} \sqrt{P(a)} + 2\sqrt{1 - P(a)}] \\ &\leq 2\sqrt{P(b)} - \sqrt{1 - P(a|b)} \sqrt{P(b)} [2\sqrt{P(a)} \\ &\quad + 2\sqrt{1 - P(a)}] \end{aligned}$$

iii) 기타의 경우는 식 (8)의 경우에서 존재할 수 없다.

따라서 i), ii), iii)에 의하여 위의 정리는 성립한다. \square

[정리 2] 만약 조건 확률 R_s 가 1이라면 R_s 의 H 값을 R_g 의 H 값을 초과할 수 없다.

$$\text{IF } P(a|b) = 1, H_s \leq H_g$$

[증명] H_g 와 H_s 의 값은 식 (6)과 식 (7)과 같이 주어진다.

$P(b) = P(ab) + P(\neg ab)$ 과 $P(a|b) = \frac{P(ab)}{P(b)} = 1$ 에 의하여 $P(\neg ab) = P(ab) - P(ab) = 0$ 이 성립한다.
따라서

$$\begin{aligned} P(a|bc) &= \frac{P(abc)}{P(bc)} = \frac{P(abc)}{P(abc) + P(\neg abc)} \\ &= \frac{P(abc)}{P(abc) + P(c|\neg ab)P(\neg ab)} = 1 \end{aligned} \quad (9)$$

식 (6)과 $P(a|b) = 1$ 에 의해서, $H_g = \sqrt{P(b)}(2 - 2\sqrt{P(a)})$

식 (7)과 식 (9)에 의하여, $H_s = \sqrt{P(bc)}(2 - 2\sqrt{P(a)})$

$P(bc) \leq P(b)$ 이므로, $H_s \leq H_g$ \square

이 법칙들은 속성 C 에 관한 추가정보없이 H_s 값의 경계를 예상하게 할 수 있게 한다. 따라서 만약 H_s 값의 경계 값이 현재의 최소 H 값(H^*) 보다 작다면 현재의 규칙에 대하여는 시스템은 더 이상 진행할 특수한 규칙들을 생성할 필요가 없다.

또한 추가로 사용되는 가지치기 방법은 현재 진행중인 규칙의 조건 확률이 1이라면, [정리 2]에 의하여 시스템은 더 이상 특수한 규칙을 생성할 필요가 없음을 알 수 있다. 가지치기 방법의 알고리즘 코드는 (그림 3)에서 설명되어있다.

```

RULES := BEST := {} ;
for each attribute  $B$  and each value  $b_i$  in attribute  $B$  do
  generate rules with  $B=b_i$  rules being left-hand side ;
  insert it into RULES;
end-for
BEST := top k rules of RULES ;
 $H^* :=$  Hellinger measure of the  $k$ th rule in BEST ;
while RULES != {} do

```

```

select a rule G from RULES ;
 $H_g :=$  Hellinger measure of rule G ;
if success rate of  $G = 1$  then delete from RULES ; /*[정리 1] */
else
    compute the bound of  $H_g$  given in Theorem 1 ;
    if  $H_g > H^*$  then
        for each attribute  $C \in$  condition part of G and each
        value  $c_i \in C$  do
            generate a specialized rule S by adding  $C = c_i$  ;
            RULES := RULES  $\cup$  S ;
            if  $H$  value of S  $> H^*$  then
                delete the bottom rule of BEST ;
                insert S into BEST ;
                update  $H^*$  ;
            end-if
        end-for
    else
        delete G from RULES ; /* [정리 2] */
    end-while
return BEST ;

```

(그림 3) 알고리즘 코드

4. 실험 결과

본 연구의 방법은 C++ 언어를 이용하여 구현되었으며 실험의 내용은 다음과 같다. 본 연구의 실험을 위하여 두 종류의 데이터를 사용하였으며 첫 번째 데이터는 (그림 4)와 같은 내용을 포함하고 있다. 이 데이터는 어느 유통업체의 실제 데이터 중에서 일부분을 발췌하여 사용하였으며 또한 제품의 숫자를 조금 축소화하였다. (그림 4)는 전체 데이터의 일부분을 보여주는데 시간 속성은 제품의 구입 시간을 연속형 숫자로 표시한 것이고 장소는 제품을 구입한 장소의 코드를 의미한다. 또한 제품 속성은 고객이 구입한 제품의 내용을 의미하며 실제 데이터는 제품의 코드를 기록하고 있다.

이와 같은 데이터를 이용하여 본 알고리즘을 수행한 결과 생성된 연속패턴 중에서 가장 상위의 10법칙을 (그림 5)에서 보여준다. (그림 5)에서 법칙 2의 경우는 제품 속성 내부에서의 연속패턴을 보여주고 있다. 즉 법칙 2의 경우는 고다치즈를 구입하는 고객은 나중에 맥주B를 구입한다는 의미이며 이는 기존의 연속패턴 생성 알고리즘에서 제공하는 기능이다. 따라서 본 알고리즘은 기존의 연속패턴 생성 알고리즘의 기능을 포함하고 있음을 알 수 있다. 나머지 법칙들은 제품 속성과 다른 속성과의 연속패턴 관계를 보여주는 법칙들이다. 예를 들어서 법칙 1의 경우는 장소 s4에서 우유를 구입한 사람은 나중에 맥주C를 구입하는 경향이 있음을 보여준다. 법칙 8의 경우에는 장소 s3에서 제품을 구입한 사람은(구입 제품의 종류에 관계없이) 나중에 맥주C를 구입하는 패턴이 있음을 보여준다. 이러한 법칙도 본 알고리즘에서 생성하는 새로운 연속패턴의 종류이다.

또한 본 알고리즘에서는 (그림 5)에서 보는바와 같이 생성되는 모든 연속패턴 법칙에 대하여 그 중요도를 H 값으로 표시하고 있다. 이와 같이 본 알고리즘은 시차관계를 가진 데이터를 읽고 좀더 일반적인 형태의 연속패턴을 자동으로 발견할 수 있음을 보여준다.

고객번호	시간	장소	제품
1	1	0	맥주B
1	1	0	사이다
1	42	0	흑맥주A
1	42	2	우유
1	77	2	독일산백포도주
1	77	3	흑맥주A
1	100	4	독일산적포도주
1	100	0	우유
1	100	0	포르투칼적포도주
...
2	2	3	맥주C
2	2	3	사이다
...

(그림 4) 연속패턴 데이터 I

No.	IF	THEN	H 값
1	장소 = s4, 제품 = 우유	제품 = 맥주C	0.0025
2	제품 = 고다치즈	제품 = 맥주B	0.0018
3	장소 = s6, 제품 = 맥주A	제품 = 흑맥주B	0.0014
4	장소 = s6, 제품 = 우유	제품 = 흑맥주B	0.0011
5	장소 = s5, 제품 = 포르투칼적포도주	제품 = 맥주B	0.0010
6	장소 = s1, 제품 = 흑맥주A	제품 = 우유	0.0009
7	장소 = s3, 제품 = 포르투칼적포도주	제품 = 맥주C	0.0008
8	장소 = s3	제품 = 맥주C	0.0008
9	장소 = s2, 제품 = 맥주B	제품 = 맥주C	0.0007
10	장소 = s6, 제품 = 흑맥주B	제품 = 맥주B	0.0007
...

(그림 5) 연속패턴 생성의 결과 I

본 연구에서 두 번째로 사용한 데이터는 제품 구매와 관련한 가상데이터를 사용하여 실험하였고 데이터는 10개의 구매와 관련한 속성으로 구성되어있으며 각 속성의 내용은 (그림 6)과 같다. 데이터 속성의 값들은 임의값을 생성하였으며 두 개의 데이터 첫 번째 데이터의 레코드 건수는 200,000 건을 사용하였으며 두 번째 데이터의 레코드 건수는 300,000 건을 사용하였다. (그림 7)은 첫 번째 데이터를 사용하여 생성한 연속패턴의 결과를 보여주며 (그림 8)은 두 번째 데이터를 사용하여 생성한 연속패턴의 결과를 보여주고있다. 앞의 실험 I 과 마찬가지로 이번 실험에서도 다양한 다차원의 연속패턴을 탐지할 수 있음을 알 수 있다.

속성	속성의 값
고객번호	A0, … A10
시간	연속값
제품	I00, I02, … I05
성별	남, 여
나이	19이하, 20대, 30대, 40대이상
결혼	기혼, 미혼
직업	공무원, 회사원, 자영업, 기타
할인여부	할인, 미할인
색상	적, 청, 황, 녹
지급수단	현금, 신용카드

(그림 6) 연속패턴 데이터 II

No.	IF	THEN	H 값
1	결혼 = 기혼	제품 = I03	0.00033
2	성별 = 남, 제품 = I03	제품 = I02	0.00031
3	직업 = 회사원	제품 = I01	0.00028
4	제품 = I01	제품 = I04	0.00026
5	할인여부 = 할인	제품 = I00	0.00020
6	색상 = 녹, 제품 = I02	제품 = I03	0.00019
7	지급수단 = 현금	제품 = I02	0.00015
8	성별 = 여	제품 = I04	0.00015
9	나이 = 20대, 제품 = I05	제품 = I00	0.00013
10	제품 = 6	제품 = I05	0.00012
…	…	…	…

(그림 7) 연속패턴 생성의 결과 II

No.	IF	THEN	H 값
1	성별 = 여	제품 = I03	0.00029
2	제품 = I01	제품 = I05	0.00027
3	색상 = 적, 제품 = I01	제품 = I04	0.00024
4	직업 = 자영업	제품 = I01	0.00024
5	지급수단 = 신용카드, 제품 = I01	제품 = I02	0.00020
6	할인여부 = 할인	제품 = I02	0.00019
7	결혼 = 기혼	제품 = I00	0.00015
8	나이 = 30대, 제품 = I00	제품 = I04	0.00012
9	제품 = I01	제품 = I01	0.00009
10	성별 = 남, 제품 = I04	제품 = I03	0.00008
…	…	…	…

(그림 8) 연속패턴 생성의 결과 III

5. 결 론

본 연구는 시차를 가지고 발생하는 데이터에 잠재하고 있는 연속패턴을 자동으로 탐지하여 법칙의 형태로 제공하는 방법을 정보이론의 헬링거 변량을 이용하여 제안하였다. 기존의 연속패턴 방법이 오직 한 개의 속성에 대해서만 연속패턴을 탐지하는데 비하여 본 연구에서 제시하는 방법은 여러 속성간의 연속패턴을 탐지하는 기능을 제공할 수 있

었으며 이는 사용자에게 훨씬 많은 정보를 제공할 수 있다. 본 연구에서 제안한 방법은 실험 데이터를 이용한 실험을 통하여 데이터속에 잠재하고 있는 연속패턴을 효과적으로 탐지할 수 있음을 알 수 있었다.

한 가지 고려할 사항으로 본 연구는 알고리즘의 수행속도를 감소시키기 위하여 두 가지 정리를 사용하여 가지치기 방법을 제공한다. 하지만 아직도 속성의 숫자가 아주 많은 데이터의 경우에는 수행속도를 더욱 감소시킬 수 있는 방법이 추가로 필요하며 수행속도의 정확한 분석과 아울러 추후의 연구과제가 될 것이다.

참 고 문 헌

- [1] Jiawei Han, Micheline Kamber, Data Mining : Concepts and Techniques, Morgan Kaufmann, August, 2000.
- [2] David J. Hand, Heikki Mannila and Padhraic Smyth, Principles of Data Mining, MIT Press, Fall, 2000.
- [3] R. Agrawal and R. Srikant, *Mining sequential pattern*, Conf. Data Engineering(ICDE '95).
- [4] R. Agrawal and R. Srikant, *Mining sequential pattern : Generalizations and Performance Improvements*, Int'l Conf. on Extending Database Technology, 1996.
- [5] R. Agrawal, R. Srikant, "Fast Algorithms for Mining Association Rules," Proc. of the 20th Int'l Conference on Very Large Databases, Santiago, Chile, Sept., 1994.
- [6] Rakesh Agrawal, Tomasz Imielinski and Arun Swami, *Mining association rules between sets of items in large databases*, In Proc. of the ACM SIGMOD Conference on Management of Data, Washington, D.C., pp.207-216, May, 1993.
- [7] C. Lee, *Learning Inductive Rules Using Hellinger Measure*, Applied Artificial Intelligence, Vol.13, No.8, pp.743-762, 1999.
- [8] R. J. Beran, Minimum Hellinger Distances for Parametric Models, *Ann. Statistics*, Vol.5, pp.445-463, 1977.
- [9] J. Han, J. Pei, B. Mortazavi-Asl, Q.Chen, U. Dayal and M.-C. Hsu., Freespan : Frequent pattern-projected sequential pattern mining, *Conf. Knowledge Discovery and Data Mining(KDD '00)*, 2000.
- [10] H. Mannila, H. Toivonen and A. I. Verkamo, *Discovery of frequent episodes in event sequences*, Data Mining and Knowledge Discovery, 1998.
- [11] M. N. Garafalakis, R. Rastogi, K. Shim, *SPIRIT : Sequential Pattern Mining with Regular Expression Constraints* Int'l COnf. on VLDB, 1999.
- [12] J. Han, J. Pei, G. Dong and K. Wang, Efficient Computation of Iceberg Cubes with Complex Measures, *Int'l Conf. on Management of Data(SIGMOD-01)*, 2001.
- [13] F. Masseglia, F. Cathala and P. Poncelet, Incremental

- Mining of Sequential Patterns in Large Databases, *European Symposium on Principles of Data Mining and Knowledge Discovery(PKDD98)*, Vol.1510, pp.176-184, 1998.
- [14] M. Zaki, N. Lesh and M. Ogihara. PLANMINE : Sequence Mining for Plan Failures, *Int'l Conf. on Knowledge Discovery and Data Mining(KDD-98)*, 1998.
- [15] M. Zaki, SPADE : An Efficient Algorithm for Mining Frequent Sequences, *Machine Learning*, Vol.42, No.1/2, pp.31-60, 2001.

이 창 환

e-mail : chlee@dgu.ac.kr

1982년 서울대학교 계산통계학과(학사)

1988년 서울대학교 계산통계학과(석사)

1994년 University of Connecticut, Dept. of Computer Science(박사)

1982년~1987년 한국기계연구소

1994년~1996년 AT&T Bell Laboratories, Middletown

1996년~현재 동국대학교 부교수

관심분야 : 기계학습, 마이닝, 생물정보학 등