

수평 분산 데이터베이스 상의 세부 데이터 유출이 없는 순차 패턴 마이닝 기법

김승우<sup>o</sup> 원정임 박상현  
연세대학교 컴퓨터과학과

{kimswo, jiwon, sanghyun}@cs.yonsei.ac.kr

Privacy Preserving Distributed Data Mining of Sequential Patterns on Horizontally Partitioned Databases

Seungwoo kim<sup>o</sup> JungIm Won, Sanghyun Park  
Department of Computer Science, Yonsei University

요약

본 논문에서는 수평 분산 데이터베이스에서 각 로컬 데이터베이스의 세부 데이터를 유출하지 않는 순차 패턴 마이닝 기법을 제안한다. 데이터 마이닝은 대용량 데이터베이스에서 유용한 지식을 추출하는 기법으로서 각 광을 받고 있다. 그러나 분산 데이터베이스를 대상으로 마이닝을 수행하는 경우, 데이터 공유에 따른 개인 혹은 집단의 프라이버시가 유출될 수 있다는 문제점이 존재한다. 따라서 본 논문에서는 프라이버시 보호를 위하여 각 로컬 데이터베이스의 세부 데이터를 보호하면서도, 마이닝 결과의 정확성을 보장할 수 있는 새로운 순차 패턴 마이닝 기법을 제안한다. 제안된 기법에서는 우선, 세부 데이터의 유출을 방지하기 위하여 마이닝의 대상이 되는 항목과 항목간의 시간 선후 관계의 성립 여부를 벡터로 표현한 후, 이들 벡터간의 스칼라 프로덕트 연산을 수행하여 얻어진 결과를 패턴의 지지도로 활용하는 방안을 제안하였다. 또한, 연산 결과에 영향을 미치지 않는 벡터를 미리 제거하여 스칼라 프로덕트 연산에 따른 비용을 감소시키는 방안을 제안하였다.

1. 서론

최근 데이터 수집 및 저장 기술의 발달로 인하여 방대한 양의 데이터들이 축적되어 저장, 관리되고 있다. 이들 대용량 데이터들을 분석하여 얻어진 정보를 사용자 의사 결정에 활용하고자 하는 연구들이 최근 활발히 진행되고 있다[1].

데이터 마이닝이란 대량의 데이터로부터 쉽게 드러나지 않는 유용한 정보들을 추출하는 과정이다[1]. 즉, 사용 가능한 모든 데이터를 기반으로 숨겨진 지식, 기대하지 못했던 패턴, 새로운 법칙과 관계를 발견하여 이를 의사 결정 과정에 활용하는 것이다. 데이터 마이닝에는 연관 규칙(association rules), 순차 패턴(sequential patterns), 군집화(clustering), 분류(classification), 예측(prediction)과 같은 여러 가지 기법들이 널리 사용되고 있다. 이 중에서 순차 패턴은 연관 규칙에 시간 개념을 포함시킨 것으로, 임의의 시간 간격을 두고 순서대로 일어난 데이터를 분석하여 발생 빈도수가 높은 패턴을 찾아내는 기법이다[2]. 실제계에서 개개인에 의해 생성되는 데이터는 임의의 시간 간격을 두고, 여러 장소에서 발생되기 때문에 이질 형태를 갖는 분산 데이터베이스에 저장되는 경우가 빈번하다. 따라서 시간의 선후 관계를 기반으로 데이터 간의 연관성 혹은 규칙성 등을 분석하는 순차 패턴을 이용하여 이들 분산 데이터베이스 내의 데이터를 통합, 분석하는 것이 유용하다.

데이터 마이닝을 통해 얻어진 정보의 신뢰성과 유용성은 분석의 대상이 되는 원본 데이터의 크기와 질에 비례한다. 따라서 최근 여러 곳에 산재되어 저장, 관리되고 있는 분산 데이터베이스에서 대량의 양질 데이터를 수집하여, 이를 대상으로 데이터 마이닝을 수행하려는 시도들이 진행되고 있다[3][4].

분산 데이터베이스에서의 데이터 마이닝 적용 사례로 금융 분야를 들 수 있다. 여러 은행들이 저장, 관리하고 있는 기존의 고객 데이터 즉, 개인 신상, 거래 실적, 대출금, 상환 일자 등의 데이터를 데이터 마이닝을 통해 분석하여 고객에 대한 자금 회수 가능성 및 신용도 등의 정보를 파악하고, 이를 이용하여 신규 대출의 가능 여부 및 대출 금액의 한도 등을 조절함으로써 대출 시에 발생할 수 있는 위험도를 낮출 수 있다. 또 다른 사례로 여러 개의 소매점을 체인망 형태로 보유하고 있는 대형 마트들을 들 수 있다. 이들 마트들은 인근의 타 마트들로의 기존 고객의 이탈을 방지하고, 신규 고객을 유치하기 위하여 여러 곳에 분산되어 있는 데이터에 대한 데이터 마이닝을 통해 고객의 구매 패턴과 선호도, 고객의 그룹별 특성 등을 분석하여 신규 서비스 개발, 패키지 상품 개발, 수익성 높은 상품 입점 등의 집중적인 마케팅 전략을 세울 수 있다.

이렇게 분산 데이터베이스에서의 데이터 마이닝은 금융, 소매, 통신, 의료 등의 다양한 분야에서 폭 넓게 활용, 적용될 수 있다. 그러나, 기업 혹은 집단의 이익, 개인의 프라이버시 보호 등의 이유로 데이터를 소유한 집단이 데이터의 공유를 제한함으로써 인하여 그 실용성에 한계가 있다. 제한적인 데이터 공유 환경에서의 마이닝 기법은 분류를 위한 의사 결정 트리 생성하는 연구에서 시작되어[5], 연관 규칙[6][7], 클러스터링[3] 등의 연구로 확장되고 있다. 그러나 순차 패턴 마이닝 기법에 대해서는 아직 미흡한 상태이며, 이와 관련된 기존 연구 [8]에서 제안된 방법은 순차 패턴 마이닝의 대상이 되는 데이터베이스

스에서 트랜잭션을 구성하는 항목이 동일하지 않은 환경, 즉 수직 분산 데이터베이스만을 대상으로 한다는 한계가 있다.

따라서, 본 연구에서는 분산 데이터베이스를 대상으로 세부 데이터의 유출을 방지하면서도, 마이닝 결과의 정확성을 보장할 수 있는 순차 패턴 마이닝 기법을 제안한다. 본 논문의 공헌은 다음과 같다. (1) 기존에 연구되지 않은 공통 항목을 갖는 수평 분산 데이터베이스를 데이터 마이닝의 대상으로 한다. (2) 데이터 공유 시, 세부 데이터의 유출을 방지하기 위한 데이터 보호 기법을 제안한다. (3) 일부의 공개된 데이터만을 이용한 순차 패턴 마이닝 기법을 제안한다. (4) 제안된 기법의 효율을 증진시키기 위한 개선 방안을 제안한다.

2. 관련 연구

분산 데이터베이스에서의 데이터 마이닝은 대량의 데이터를 대상으로 마이닝을 수행하여 보다 유용하고, 신뢰성 있는 정보를 추출하기 위하여 연구되어져 왔으며, 일반적으로 두 가지 방법이 있다. 첫 번째는 분산되어 있는 각 로컬 데이터베이스내의 모든 데이터들을 중앙의 한 데이터베이스에 저장해 놓은 후, 이를 데이터 마이닝하는 것이다[9]. 이 방법은 구현이 쉽고, 간단하지만 대량의 데이터를 중앙으로 전송하기 위한 통신 비용 및 저장 공간의 오버헤드가 크다는 단점이 있다. 두 번째는 각각의 로컬 데이터베이스를 대상으로 마이닝을 수행한 후, 그 결과를 집계하여 결과를 얻는 것이다[10]. 이 방법은 통신 비용 및 저장 공간에서 효율적이나, 데이터 마이닝을 통해 얻어진 결과의 정확성이 떨어진다라는 단점이 있다.

데이터 마이닝에 있어서의 프라이버시 문제는 두 가지 측면에서 제기되었다. 첫 번째는 마이닝의 대상이 되는 데이터에 의해 개인의 사생활 정보가 유출되어 상업적인 용도로 도용되거나, 악의적인 사용자에게 의하여 악용되는 것을 방지하는 것이며[11], 두 번째는 다수의 기업 혹은 집단의 데이터를 마이닝하는 경우, 각 집단의 비공개 데이터가 다른 집단으로 유출되는 것을 방지하는 것이다[5]. 분산 데이터베이스를 대상으로 각 로컬 데이터베이스내에 저장되어 있는 데이터들의 프라이버시를 보장하면서, 이들을 마이닝하기 위한 연구로 참고 문헌 [6][7]에서는 연관 규칙을 이용한 마이닝 기법을 제안하였으며, 참고 문헌 [4]에서는 나이브 베이즈 분류(Naive Bayes classification)를 이용한 마이닝 기법을 제안하였다. 또한, 참고 문헌 [3]에서는 군집화를 위한 마이닝 기법을 제안하였다.

특히, 참고 문헌 [6]에서는 로컬 데이터베이스 간에 공통 항목이 존재하는 수평 분산 데이터베이스를 대상으로 연관 규칙 마이닝 기법을 제안하였으며, 이를 위해 마이닝의 대상이 되는 각 데이터의 출처를 공개하지 않으면서도, 전체 데이터를 마이닝할 수 있는 상호 암호화(commutative encryption) 기법과 각각의 로컬 데이터베이스로부터 얻어진 지지도로부터 전체 지지도를 계산할 수 있는 시큐어섬(secure sum) 기법을 제안하였다. 또한 참고 문헌 [7]에서는 공통 항목이 존재하지 않는 수직 분산 데이터베이스를 대상으로 하는 연관 규칙 마이닝 기법을 제안하였으며, 이를 위해 데이터베이스 내에서 각 항목(item)의 발생 여부를 벡터 형태로 표현한 후, 스칼라 프로덕트(scalar product) 연산을 수행하여 지지도를 계산하는 방법과

연산 시에 벡터 값이 공개되는 것을 방지하기 위한 보안 프로토콜을 제안하였다. 또한, 참고 문헌 [8]에서는 분산 데이터베이스를 위한 순차 패턴 마이닝 기법을 제안하였다. 각 로컬 데이터베이스 내에서 트랜잭션의 대상이 되는 데이터를 병합하여 시간 순서에 따라 오름차순으로 정렬한 다음, 정렬된 데이터를 발생 시간과 발생 여부로 분리하여 벡터 형태로 표현한 후, 스칼라 프로덕트 연산에 의하여 지지도를 계산하는 방법을 제안하였다. 이 때, 데이터의 유출을 방지하기 위하여 데이터를 매핑 테이블을 이용하여 일련의 숫자로 변환하였으며, 랜덤 방식에 의해 생성된 임의의 값을 데이터의 병합 과정에 추가하는 방식을 사용하였다. 그러나, 제안된 방법은 로컬 데이터베이스 간에 공통 항목이 존재하지 않는 경우만을 대상으로 하며, 마이닝의 대상이 되는 데이터를 어느 로컬 데이터베이스에서 저장하고 있는지를 알 수 있으며, 마이닝에 의하여 생성된 패턴에 대하여 역 매핑을 수행하는 경우에 데이터의 유출을 방지하기 위하여 사용된 매핑 값을 유추해 낼 수 있으므로 프라이버시를 보장할 수 없다는 문제점이 있다.

본 논문에서는 기존에 연구되지 않은 공통 항목을 갖는 수평 분산 데이터베이스에서의 순차 패턴 마이닝 기법을 제안한다. 제안된 기법에서는 프라이버시를 보호하기 위하여는 참고 문헌 [6]에서 제안된 상호 암호화 기법과 시큐어 섬 기법을 사용하여 데이터의 출처를 감추고, 순차 패턴 마이닝에서의 지지도 계산을 위하여는 수직 분산 데이터베이스를 대상으로 제안된 [7]의 알고리즘을 확장, 적용한다. 수평 분산 데이터베이스의 경우, 각 로컬 데이터베이스내에 같은 값을 가지는 공통 항목이 출현할 수 있으므로, 참고 문헌 [8]에서와 같이 한 번의 스칼라 프로덕트 연산으로는 지지도를 구할 수 없다. 따라서, 본 연구에서는 스칼라 프로덕트 연산을 위한 알고리즘을 확장, 적용하는 방안을 제안하였으며, 연산을 위하여 생성된 벡터에서 연산 결과에 영향을 미치지 않는 부분을 미리 제거하여 연산 비용을 감소시키는 방안을 제안하였다.

3. 제안하는 기법

본 논문에서는 각 로컬 데이터베이스간에 공통 항목이 존재하는 수평 분산 데이터베이스를 대상으로 하며, 순차 패턴 마이닝 기법을 이용하여 로컬 데이터베이스내의 데이터들간의 관계, 규칙 등의 의미있는 유용한 정보를 찾아내고자 한다. 이를 위해 여러 개의 소매점을 체인망 형태로 보유하고 있는 대형 마트들을 예로 든다. 각 마트들은 <표 1>과 같이 고객 ID, 구매 시간, 구매 품목 등으로 구성된 구매 정보를 관리한다. 이때, 각 마트들간에는 같은 구매 품목들 즉, 공통 항목들이 존재한다. 이를 이용하여 고객의 구매 패턴 분석을 위한 마이닝을 수행할 경우, 지지도 계산을 위하여 고객 ID, 구매 품목 등은 반드시 공개되어야 한다. 그러나 각 마트의 수익을 위하여 해당 데이터가 공개되는 것을 제한할 수 있다. 따라서 본 논문에서는 제한적인 데이터 공유만을 허락하는 이들 데이터의 유출을 방지하면서도, 마이닝 결과의 정확성을 보장할 수 있는 순차 패턴 마이닝 기법을 제안한다.

3.1 수직 분산 데이터베이스를 위한 기법

참고 문헌 [8]에서는 각 로컬 데이터베이스간에 공통 항목이 존재하지 않는 수직 분산 데이터베이스를 대상으로 하는 순차 패턴 마이닝 기법을 제안하였다. 제안된 기법에서는 우선, 상제 데이터의 유출을 방지하기 위하여 각 로컬 데이터베이스내의 모든 항목을 매핑 테이블을 이용하여 유일한 값으로 매핑한다. 예를 들어 <표 1>의 마트 1에서의 {전자사전}은 1A, 마트 2의 {전기밥솥, 전자사전}은 2B, 마트 3의 {건전지}는 2C 등의 유일한 값으로 매핑한다. 다음, Apriori 알고리즘 [2]에 이용하여 매핑 값을 기반으로 후보 패턴 생성하고, 생성된 후보 패턴에 대한 지지도를 계산하여, 최소 지지도 이상을 만족하는 패턴만을 빈번 패턴으로 생성한다.

지지도 계산을 위하여 각 로컬 데이터베이스내의 모든 트랜잭션에 대한 발생 시간을 공유하며, 이를 이용하여 각 패턴의 실제 발생 유무를 표현하는 벡터와 패턴간의 시간 선후 관계가 성립하면 1, 성립되지 않으면 0의 값을 갖는 벡터를 생성한다. 생성된 벡터간의 스칼라 프로덕트 연산을 수행하여 얻어진 값이 해당 패턴의 지지도가 된다. 예를 들어, 매핑 값으로 변환된 패턴 {2A->2C}에 대한 지지도는 패턴 2A의 발생 시간 (20일 14시 25분, 20일 14시 30분, 26일 13시 40분)와 패턴 2C의 발생 시간 (21일 16시 25분, 27일 11시 30분, 24일 17시 20분)에 대하여 고객 ID별 시간 값을 비교하여 2A의 발생 시간이 2C의 발생 시간 보다 이전인 경우 1, 아닌 경우 0으로 표현한 벡터 (1, 1, 0)를 구한다. 구해진 시간 벡터 (1, 1, 0)를 패턴 2A의 발생 여부를 표현한 벡터 (1, 1, 0)과 패턴 2C의 발생

여부를 표현한 벡터 (0, 1, 0)와 함께 스칼라 프로덕트 연산을 수행하여 패턴 {2A->2C}의 지지도로 1의 값을 얻는다.

<표 1> 각 마트의 고객 구매 정보

마트 1		
고객 ID	구매 시간	구매 품목
1	2월 20일 14시 25분	전기밥솥, 전자사전
2	2월 20일 14시 30분	전자사전
1	2월 26일 12시 20분	건전지, 쌀
마트 2		
고객 ID	구매 시간	구매 품목
3	2월 21일 16시 25분	전기밥솥, 전자사전
3	2월 23일 14시 25분	건전지
마트 3		
고객 ID	구매 시간	구매 품목
2	2월 22일 15시 30분	전기밥솥
2	2월 27일 11시 30분	건전지

3.2 수평 분산 데이터베이스를 위한 기법

3.1 절에서 제시된 기법을 본 논문에서 다루는 수평 분산 데이터베이스에 적용하는 경우 다음과 같은 문제점들이 발생한다. (1) 서로 다른 로컬 데이터베이스내에 동일 항목이 존재하는 경우, 이들 항목들은 매핑 테이블을 통해 서로 다른 매핑 값으로 변환된다. 따라서, 동일 항목이 서로 다르게 표현되므로 마이닝 결과의 정확성을 보장할 수 없다. (2) 임의의 항목에 대한 지지도 계산을 위하여 수직 분산 데이터베이스에서는 항목을 저장하고 있는 단 하나의 로컬 데이터베이스만을 액세스하여, 항목에 대한 발생 빈도를 구하면 된다. 그러나, 수평 분산 데이터베이스에서는 임의의 항목이 여러 로컬 데이터베이스에 저장되어 있을 수 있으므로, 다수의 로컬 데이터베이스를 액세스하여 항목에 대한 전체 발생 빈도를 구해야한다. 즉, 더 많은 데이터 전송 비용과 스칼라 프로덕트 연산에 대한 비용이 요구된다. (3) 매핑 값을 기반으로 생성된 빈번 패턴에 대하여 매핑 테이블을 이용한 역 매핑을 수행하는 경우, 실제 항목에 대한 정보가 유출될 수 있다.

<표 2> 마트 1, 2, 3의 구매 여부에 따른 표현

마트1									
고객 ID	전기밥솥		전자사전		쌀		건전지		
	1	1	2월 20일 14시 25분	1	2월 20일 14시 25분	1	2월 26일 12시 20분	1	2월 26일 12시 20분
2	0	2월 24일 11시 30분 <sup>1)</sup>	1	2월 20일 14시 30분	0	2월 26일 13시 20분 <sup>1)</sup>	0	2월 24일 14시 50분 <sup>1)</sup>	
3	0	2월 22일 16시 25분 <sup>1)</sup>	0	2월 26일 13시 40분 <sup>1)</sup>	0	2월 21일 12시 40분 <sup>1)</sup>	0	2월 23일 11시 10분 <sup>1)</sup>	
마트2									
고객 ID	전기밥솥		전자사전		쌀		건전지		
	1	0	2월 22일 19시 35분 <sup>1)</sup>	0	2월 24일 17시 45분 <sup>1)</sup>	0	2월 20일 13시 20분 <sup>1)</sup>	0	2월 24일 15시 40분 <sup>1)</sup>
2	0	2월 24일 17시 40분 <sup>1)</sup>	0	2월 26일 12시 00분 <sup>1)</sup>	0	2월 26일 18시 10분 <sup>1)</sup>	0	2월 23일 12시 20분 <sup>1)</sup>	
3	1	2월 21일 16시 25분	1	2월 21일 16시 25분	0	2월 23일 14시 25분 <sup>1)</sup>	1	2월 23일 14시 25분	
마트3									
고객 ID	전기밥솥		전자사전		쌀		건전지		
	1	0	2월 27일 17시 45분 <sup>1)</sup>	0	2월 25일 19시 20분 <sup>1)</sup>	0	2월 22일 18시 40분 <sup>1)</sup>	0	2월 21일 16시 25분 <sup>1)</sup>
2	1	2월 22일 15시 30분	0	2월 23일 15시 30분 <sup>1)</sup>	0	2월 23일 15시 30분 <sup>1)</sup>	1	2월 27일 11시 30분	
3	0	2월 24일 17시 40분 <sup>1)</sup>	0	2월 23일 17시 40분 <sup>1)</sup>	0	2월 26일 18시 30분 <sup>1)</sup>	0	2월 24일 17시 20분 <sup>1)</sup>	

이러한 문제점 해결을 위하여 본 논문에서는 기존 [8]의 방식을 확장, 적용한 새로운 기법을 제안한다. 이를 위해 우선, 길이 1인 빈번 패턴 생성을 위하여 참고 문헌 [6]에서의 상호 암호화(commutative encryption) 기법과 시큐어 섬 기법을 사

1) 임의로 생성된 시간 데이터

