

디자인 자료의 분석 도구로서의 데이터마이닝 기법의 적합성 점검*

이태경
디자인학부 디지털정보디자인전공

<요약>

특정 상품이 시장에서 타 상품과의 경쟁에서 이기기 위해서는 소비자 기호를 충족시키는 디자인이 중요한 요소이다. 소비자 기호를 분석하기 위해서 사용되는 회귀분석을 하기 위해서는 회귀분석모델을 설정하여야 한다. 회귀분석모델을 설정하기 위해서는 분석자료의 영역에 대한 지식과 회귀분석에 대한 지식이 필요하다. 이는 일반 분석자에게는 큰 제약이 된다. 이 논문에서는 데이터마이닝 기법을 이용하여 회귀분석모델을 쉽게 할 수 있는 방법을 시도한다.

The suitability of a datamining method for the analysis tool of design-related data

Lee, Tae-Kyong
Dept. of Digital Information Design, University of Ulsan

<Abstract>

A design is a key factor of a product for the product to dominate against other products in a market. A regression methodology is used to analyse a consumer's trend in a market. A analyzer's knowledge of both the domain of a data and the regression methodology itself is the least requirement for a reasonable result of a regression. This may be a burden and sometimes hinder a layman from trying a regression. This paper

* 이 논문은 2001년 울산대학교의 연구비에 의하여 연구되었음

proposes a method to integrate a datamining technique into a regression methodology and to use a regression methodology with easy.

1. 서론

최근 특정 산업 분야에서 두각을 나타내는 기업들의 제품을 살펴보면 그 상품 자체의 기능성의 우수성뿐만 아니라 소비자들의 기호에 부응하는 디자인의 우수성도 뛰어나다는 것을 알 수 있다. 이는 특정 상품이 시장에서 타 상품과의 경쟁에서 이기기 위해서는 소비자 기호를 충족시키는 디자인의 개발이 필요함을 의미한다. 특히 시장이 개방화되어 무한 경쟁 시대인 요즈음 국가간의 제품 생산의 기술 격차가 좁혀지고 있는 현재 다양한 소비자 요구에 부응하여 상품을 다양화, 차별화 시킬 수 있는 디자인 개발의 중요성은 더욱더 커지고 있다. [5]

다양한 소비자 기호에 부응하는 상품 생산을 하기 위해서는 디자인에 대한 다양한 소비자들의 욕구를 정교하게 분석할 수 있는 방법론이 제공되어야 한다. 이 방법론은 디자인 분야의 자료들이 가지는 특성을 처리할 수 있어야 한다. 그리고 이 자료들을 가지고 실지로 분석을 행하는 디자인 분야의 종사자들이 그 방법론을 쉽게 사용할 수 있어야 한다.

다양한 소비자 기호를 정확히 분석하기 위해서는 가능한 다양한 자료를 가지고 직관적인 분석 방법을 지향하고 체계적인(systematic) 분석 방법과 기존의 밝혀진 자료들과의 인과관계 이외의 관계를 밝힐 수 있는 방법을 사용하여야 한다. 이 방법에 의한 분석이라 함은 직관적인 관찰과 기존의 인과관계에 의해서 알 수 있는 분석을 뛰어넘어 자료에는 내포되어 있으나 관찰되고 발견될 수 없었던 관계까지의 분석을 의미한다.

자료간의 상관관계의 분석에 가장 많이 사용되는 방법론은 회귀분석(regression)이다.[1,2] 이 논문이 다루는 응용분야에 적용하면 소비자 선호(output)와 그 선호를 결정하는 요소(input)들간의 상관관계를 찾는 것이다. 그리고 그 상관관계를 나타내는 함수를 이용하여 미래의 추세를 예측하는 것이다.

일반적으로 회귀분석의 성공여부는 설정된 회귀분석의 관계식이 실지의 자료간의 관계를 얼마나 정확하게 반영하는 가에 달려있다. 이는 자료 자체에 대한 전문가적인 지식과 회귀분석 자체에 대한 지식과 경험이 회귀분석의 성공 여부를 결정한다는 것이다. 이는 디자인 자체에 관한 지식과 이 지식을 회귀분석의 모델로 표현할 수 있는 지식이 필요함을 의미한다. 왜냐하면 회귀분석에서는 전문가가 파악하지 못하는 자료간의 관계는 모델에서 제외되고 이는 그 자료들간의 관계는 분석의 대상에서 제외되기 때문이다.

데이터마이닝의 여러 분야 중에서 여러 요소들간의 관계(association)를 찾는 분야가 있다.[3,4] 데이터마이닝에서 연관관계를 찾는 방법들은 주어진 자료에만 근거하여 여러 자료들 간의 관계를 찾기 때문에 기존에 알려진 자료들간의 관계에 근거한 소비자 분석이 간과하기 쉬운 관계도 찾아낸다. 이는 소비자 분석의 관점에서 보면 분석 결과의 유의성은 분석자의 지식과 경험을 포함은 어떠한 사전 편견에도 영향을 받지 않는다는 것을 의미한다. 즉, 데이터마이닝 방법을 이용하면 새로운 관계를 찾을 수 있는 소비자 분석이 가능하다는 것을 의미한다.

회귀분석은 정확한 자료들간의 관계를 도출할 수 있으나 먼저 모델 설정이 정확하여야 한다. 이 모델 설정은 전문가의 지식과 경험에 의존한다. 데이터마이닝 기법은 분석가의

어떠한 사전 편견에도 의존하지 않고 자료 자체에 존재하는 관계를 찾아낸다. 그러나 정확한 관계의 도출은 불가능하다. 이에 이 논문은 데이터마이닝 기법을 이용하여 먼저 자료들 간의 관계를 찾고 이 정보를 이용하여 회귀분석의 모델설정을 하는 방법을 제시한다.

2장에서는 일반적인 회귀분석 모델을 간단히 설명하며 왜 모델 설정이 어려운지 설명한다. 3장에서는 데이터마이닝 기법의 연관규칙(association rule)을 찾는 알고리즘을 설명하고 이를 회귀분석 모델의 설정에 사용하는 방법을 제시한다. 그리고 4장은 결론이다.

2. 회귀분석의 모델

일반적인 회귀분석의 모델은 $y_t = a_{11}x_{11} + \dots + a_{1n}x_{1n} + a_{21}x_{21} + \dots + a_{2m}x_{2m} + a_{31}y_{t-1} + \dots + a_{3z}y_{t-z}$ 로 표현된다. 이 표현에서 a_{21}, \dots, a_{2m} 중 하나라도 0이 아니면 비선형(non-linear) 모델이 되고 a_{31}, \dots, a_{3z} 중 하나라도 0이 아니면 종속변수 y 의 타임래그가 있는 모형이 된다.[1,2] 이 논문에서는 회귀분석에서 가장 간단한 모델인 선형 비타임래그의 모델을 이용한다.

선형 비타임래그의 모델은 $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n$ 으로 표현된다. 이 모델에서 어떤 상품의 만족도를 나타내는 변수는 종속변수 Y 이고 만족도를 결정하는 변수는 X_1, \dots, X_n 이라고 하자. 회귀분석은 주로 다음과 같은 목적을 위하여 사용된다.

- (1) Y 의 변이에 대한 최상의 예측을 할 수 있도록 X 들에 대한 Y 의 회귀관계를 수립하는 것.
- (2) Y 에 영향을 끼치는 독립변수 X 가 여러 가지로 그 수가 대단히 많을 때, 실은 그들 중에서 적절히 선택하면 몇 개 안되는 독립변수 X 들만으로 충분한 경우가 있다. 독립변수 50개 이상이 고려대상으로 될 경우도 불과 3, 4개의 독립변수 X 에 관한 회귀 Y 를 생각하는 것으로 충분한 경우가 있다는 말이다.
- (3) 다중회귀를 고찰하는 목적이 예측을 위한 것이 아닌 경우도 있다. 어떤 독립변수 X 가 Y 와 가장 관계가 깊은가 하는 것을 찾아내고 그 중요성의 순위를 정하려는 것이 목적인 경우도 있다.

회귀방정식을 찾고자 하는 자료를 가지고 얻은 표본회귀방정식이 $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_n X_n$ 이라하자. 이 표본회귀방정식은 회귀분석의 목적 (1), (2), (3)을 만족하여야 한다. 표본회귀방정식이 표본자료를 얼마나 정확히 설명하고 미래의 예측력을 가지고 있는가를 나타내는 척도는 표본상관계수이다. 이는 회귀분석 목적 (1)에 해당된다. 그리고 회귀 목적 (2), (3)은 회귀모델의 추정치 $\hat{\beta}_i$ 의 절대치와 관계가 있다.

표본회귀방정식이 실지의 (과거의) 자료를 얼마나 정확하게 설명하고 있는가는 실지의 자료 Y 의 값과 \hat{Y} 의 차이에 의해서 결정된다. 이것은 표본상관계수($-1 \leq r \leq 1$)로 표현된다. 표본상관계수는 직선회귀방정식에서 두 변량 X_i 와 Y 의 상관관계의 척도이기 때문이다. 표본상관관계의 절대값이 클수록 X_i 와 Y 의 관계를 밀접하게 표본회귀방정식이 설명하는 것이며 따라서 표본상관계수는 회귀분석의 목적 (1)에 관련이 있다.

표본회귀방정식에서 $\hat{\beta}_i$ 의 크기는 X_i 의 변화가 Y 의 값에 얼마나 영향을 미치는 가를 나타낸다. 이는 회귀분석의 목적 (2)와 (3)과 관련이 있다. 만약, 회귀모델을 설정하기 전에

고려의 대상이 되는 X_i 중에서 중요도의 순서를 알 수 있다면 모든 X_i 를 사용하지 않고 간단하게 목적 (2)와 (3)을 충족시킬 수 있다.

회귀분석은 1)모델을 설정하고, 2)분석대상이 되는 자료에 관하여 모델의 설명력을 높이기 위해 수 차례의 조정(fine-tuning)을 통하여 그 모델에 근거한 회귀분석함수를 구한다. 이는 앞에서 설정된 회귀분석의 목적 (2)와 (3)의 관점에서 어떤 X_i 를 회귀모델에 넣어서 목적 (1)을 만족시키는 표본회귀방정식을 구하는 것이다. 이는 표본회귀방정식을 구할 때 X_i 의 선택이 표본회귀방정식의 결과에 영향을 미친다는 것을 말한다.

다음 3장에서는 데이터마이닝의 연관관계와 연관관계추출을 위한 알고리즘을 설명한다. 그리고 이 연관관계에서 얻어지는 결과를 어떻게 회귀모델의 설정할 때 사용할 수 있는지 설명한다.

3. 연관관계(Association Rule) 추출을 위한 알고리즘.

항목들 사이의 연관관계라 함은 어떤 사건에서 특정항목이 나타날 때 다른 어떤 항목이 같이 나타남을 의미한다. 데이터마이닝의 연관관계 추출에 관해서 설명할 때 가장 빈번히 사용되는 소비행태에 관한 예를 가지고 설명한다. 상점에서 판매되는 상품의 집합 $I=\{i_1, i_2, \dots, i_m\}$ 이다. 그러면 어떤 하나의 소비는 $(t_i, \text{ 이후 } “트랜잭션”)$ 이라는 용어로 대체 사용된다.) I 의 부분집합으로 표시될 수 있다. 그리고 모든 소비(T)는 $T=\{t_1, \dots, t_n\}$ 으로 표현될 수 있다. 그리고 각각의 t_i 에는 고유의 번호(TID)가 부여된다.

정의1: X 는 I 의 부분집합이며 k 는 X 가 t_i 의 부분집합인 t_i 의 수라 하자. $\text{supp}(X)$ 는 k/n 으로 정의된다.

즉, $\text{supp}(X)$ (지지도)는 모든 소비에서 어떤 상품의 집합 X 를 포함하는 비율을 나타낸다. 데이터마이닝에서는 사용자가 정하는 최소지지도(Smin) 이상의 지지도를 가지는 상품집합 X 만을 고려한다. 이유는 관심있을 정도로 빈발하게 나타나는 상품함집만을 고려하기 때문이다.

어떤 소비에서 X 를 소비하면 Y 를 소비하는 행태는 연관규칙 $R: X \Rightarrow Y$ (X 이면 Y 이다.)로 표현된다. 예를 들어 X 가 {땅콩, 오징어} 이면 Y 는 {맥주, 소주, 맥실주}가 될 수 있을 것이다. 이때 위의 규칙이 소비의 집합 T 에서 중요성을 가지기 위해서는 어떤 소비가 X 를 가지고 있을 때 또한 Y 를 가지는 비율이 얼마나가 하는 것이다. 이 개념은 규칙의 신뢰도(confidence)개념으로 표현된다.

정의2: 만약 소비의 집합 T 에서 X 가 나타날 때 Y 가 나타나는 확률을 연관규칙 $R: X \Rightarrow Y$ 의 신뢰도이며 $\text{conf}(R)$ 로 나타낸다.

이해를 높이기 위하여 $\text{supp}(X) = (\text{number of transactions such that } X \subseteq I) / (\text{number of transactions})$ 이며 $\text{conf}(R) = p(Y \subseteq T | X \subseteq T)$ 이다. 즉, $\text{conf}(R)=p(Y \subseteq T | X \subseteq T) = P(Y \subseteq T \wedge X \subseteq T) / P(X \subseteq T) = \text{supp}(X \cup Y) / \text{supp}(X)$ 이다.

연관규칙의 추출은 다음의 두 단계로 성취된다.

1. 빈발 항목집단(large item sets)을 찾아낸다. 항목들의 전체집합 I 의 부분집합이면서 몇 개의 항목들로 구성된 항목들의 모임을 항목집합이라 무른다. 미리 결정된 최소 지지도 (Smin) 이상의 지지도를 가지는 항목집합들의 모든 집합들을 빈발 항목집합이라 한다.

2. 모든 빈발 항목집합 l에 대해서 l의 공집합이 아닌 부분집합들을 찾는다. 각각의 그러한 부분집합 a에 대하여 만약 supp(a)의 비율이 적어도 최소 신뢰도 Cmin 이상이면, 즉, $\text{supp}(l)/\text{supp}(a) \geq C_{\min}$, $a \Rightarrow (l - a)$ 의 형태의 규칙을 출력한다.

2.1 모든 빈발 항목집합 찾기

잠재적인 빈발 항목집합들의 수는 모든 항목들의 면집합(power set)의 크기와 같다. 이 빈발 항목집합을 찾는 알고리즘들은 후보(candidate)라 불리우는 빈발 가능성이 있는 항목집합들을 생성하고 그 항목 중에서 최소지지도를 가지는 항목들을 선택한다. 이 논문에서는 Apriori[3,4] 알고리즘을 이용하여 빈발 항목집합을 찾는다.

데이터베이스안에 그림과 같은 소비집합이 저장되어 있다고 하자. 그리고 모든 상품집합 $I=\{A, B, C, D, E\}$ 이다.

TID	Items
100	A C D
200	B C E
300	A B C E
400	B E

그림1 트랜잭션 데이터베이스

C1																										
Scan Database --->	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{A}</td> <td>2</td> </tr> <tr> <td>{B}</td> <td>3</td> </tr> <tr> <td>{C}</td> <td>3</td> </tr> <tr> <td>{D}</td> <td>1</td> </tr> <tr> <td>{E}</td> <td>3</td> </tr> </tbody> </table>	Item Set	Support	{A}	2	{B}	3	{C}	3	{D}	1	{E}	3	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{A}</td> <td>2</td> </tr> <tr> <td>{B}</td> <td>3</td> </tr> <tr> <td>{C}</td> <td>3</td> </tr> <tr> <td>{E}</td> <td>3</td> </tr> </tbody> </table>	Item Set	Support	{A}	2	{B}	3	{C}	3	{E}	3		
Item Set	Support																									
{A}	2																									
{B}	3																									
{C}	3																									
{D}	1																									
{E}	3																									
Item Set	Support																									
{A}	2																									
{B}	3																									
{C}	3																									
{E}	3																									
C2																										
Scan Database --->	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{A B}</td> <td>1</td> </tr> <tr> <td>{A C}</td> <td>2</td> </tr> <tr> <td>{A E}</td> <td>1</td> </tr> <tr> <td>{B C}</td> <td>2</td> </tr> <tr> <td>{B E}</td> <td>3</td> </tr> <tr> <td>{C E}</td> <td>2</td> </tr> </tbody> </table>	Item Set	Support	{A B}	1	{A C}	2	{A E}	1	{B C}	2	{B E}	3	{C E}	2	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{A C}</td> <td>2</td> </tr> <tr> <td>{B C}</td> <td>2</td> </tr> <tr> <td>{B E}</td> <td>3</td> </tr> <tr> <td>{C E}</td> <td>2</td> </tr> </tbody> </table>	Item Set	Support	{A C}	2	{B C}	2	{B E}	3	{C E}	2
Item Set	Support																									
{A B}	1																									
{A C}	2																									
{A E}	1																									
{B C}	2																									
{B E}	3																									
{C E}	2																									
Item Set	Support																									
{A C}	2																									
{B C}	2																									
{B E}	3																									
{C E}	2																									
C3																										
Scan Database --->	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{B C E}</td> <td>2</td> </tr> </tbody> </table>	Item Set	Support	{B C E}	2	<table border="1" style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th>Item Set</th> <th>Support</th> </tr> </thead> <tbody> <tr> <td>{B C E}</td> <td>2</td> </tr> </tbody> </table>	Item Set	Support	{B C E}	2																
Item Set	Support																									
{B C E}	2																									
Item Set	Support																									
{B C E}	2																									

그림2 후보 항목집합의 생성과 빈발항목집합

한 패스(pass)에서 발견된 빈발항목집합들이 다음 패스를 위한 후보 집합을 생성하는 기

반으로 사용되고 이러한 과정은 반복적으로 수행된다. 위의 예에서 처음 패스의 후보집합(C1)에서 최소지지도 2보다 큰 항목들이 처음 패스의 빈발항목집합(L1)이 된다. L1을 바탕으로 두 번째 패스의 후보집합(C2)가 만들어지며 C2의 각각의 항목의 지지도를 검사하고 최소지지도 2보다 큰 항목들이 두 번째 패스의 빈발항목집합(L2)이 된다. 같은 방법으로 세 번째 빈발항목집단(L3)이 구해진다. L3에 있는 항목이 하나이기 때문에 더 큰 빈발항목집단의 가능성은 없다. 그렇기 때문에 알고리즘은 이 순간에 멈춘다.

3. 빈발항목집단과 선형회귀모델

이 논문에서 고려하고 있는 회귀모델은 $Y = \beta_0 + \beta_1X_1 + \dots + \beta_nX_n$ 인 선형다중회귀(multiple regression)이다. 데이터마이닝 기법을 이용하여 빈발집단항목속에 포함되는 항목을 회귀모델에서 사용되는 독립변수 X로 선택하면 회귀모델의 대상이 되는 자료에 관한 전문적인 지식이 없어도 회귀모델의 설정이 용이하게 될 가능성이 높다.

실제의 자료를 가지고 데이터마이닝의 연관규칙을 찾는 알고리즘의 처음단계에서 빈발항목집단에 포함되는 항목이 회귀모델의 독립변수 X의 설정에 도움이 되는지 살펴보자. [6]에 인터넷 디자인 요소가 구매결정에 미치는 영향을 알아보기 위하여 설문조사를 회귀분석한 결과가 있다. 설문지는 다음과 같다.

설문의 내용은 전자상거래에서 전자상거래의 구매결정에 웹의 디자인이 얼마나 영향을 미치는가에 관한 질문내용이다. 설문 응답자는 구매결정에 영향을 미치는 각 항목에 표시를 하며 선택된 칸에 대해서 각각 1에서 5까지의 값이 부여된다.

이 자료와 연관규칙을 설명하기 위하여 사용된 소비행태에 관한 자료와는 차이점이 있다. 설문자료에서 각 항목은 구매자료의 각각의 상품에 해당된다고 볼 수 있다. 그러나 소비행태 자료에서는 각 상품의 구매와 불구매 결정에 따라서 두 가지의 선호만 나타난다. 그러나 설문 자료에서는 분석요구에 따라서 선호값을 여러 가지로 나눌수 있다. 따라서 설문자료의 선호값 중에서 분석자가 정한 값(threshold)을 상회하는 항목만 선택되었다고 생각하고 연관규칙 알고리즘을 적용한다.

이 자료를 가지고 마지막으로 얻은 빈발항목집단에 남는 항목은 설문 응답자가 생각하는 웹의 디자인이 전자상거래에서 구매에 영향을 미치는 주요항목들의 집합으로 생각될 수 있다. 그리고 이 항목들을 가지고 회귀모델을 설정하여 표본회귀방정식을 구하면 회귀를 하는 목적 (1), (2), 그리고 (3)을 성취할 수 있다.

	전혀 그렇지 않다		보통이다.		매우 그렇다.
	1	2	3	4	5
네비게이션도움					
그래픽					
색상					
멀티미디어					
레이아웃					
홈페이지 구조					

그림3 전자상거래 구매결정 설문지

위의 설문 결과를 알고리즘 Apriori[3,4]에 적용할 때 항목의 선택, 불선택을 위하여 threshold를 사용하였다. 적은 값을 사용하면 마지막에 남는 빈발항목집단에 속하는 항목의 수가 많으며 반대로 큰 값을 사용하면 빈발항목집단에 속하는 항목의 수가 많다. 이는 회귀모델의 설정에서 독립변수 X의 수에 영향을 미친다. 그러나 여전히 분석자가 설정한 기준하에서 Y에 영향을 가장 많이 미치는 독립변수의 집합이다.

4. 구현과 결과

이 논문에서 검증하여야 하는 것은 데이터마이닝 연관규칙을 찾기 위해서 사용되는 빈발항목집단의 알고리즘을 이용하여 얻어지는 빈발항목집단에 속하는 항목들이 실제로 회귀분석의 모델에 이용할 수 있는가 하는 것이다. 이는 다음과 같은 방법으로 검증한다.

- (1) 설문 내용에 있는 모든 항목들을 독립변수 X로 설정하여 표본회귀방정식을 구한다.
- (2) threshold값을 변화시켜 여러 가지의 값에 대해서 빈발항목집단 추출작업을 수행한다.
- (3) 각각의 threshold 값에 따라서 구해지는 빈발항목집단에 속하는 항목에 대응하는 독립변수 X_i 가 있다. 만약, 빈발항목집단에 속하는 항목이 n개이면 이 항목들에 대응하는 표본회귀방정식의 독립변수 X_i 의 β_i 의 값은 상위 n번째의 값이다.
- (3)을 주장하는 이유는 다음과 같다. 높은 threshold를 설정하면 설문지에서 설문 응답자들이 강력하게 의견을 피력하는 항목만 빈발항목집단에 선택된다. 이를 표본회귀방정식을 가지고 해석하면 그 항목에 해당하는 X_i 의 값이 빈발항목집단에 선택되지 않은 항목에 해당하는 X_j 의 값보다 Y의 값에 더 큰 영향을 끼쳐야 한다는 것이다. 이는 β_i^* 의 값이 β_j^* 의 값보다 커야 한다는 것이다.

이 주장을 검증하기 위하여 실제 자료를 가지고 표본회귀방정식을 구해야 하며 빈발항목집단을 추출하기 위하여 Apriori[3,4]알고리즘을 구현하여야 한다. 실지 자료는 3장에서 언급된 웹디자인이 전자상거래 구매 결정에 관한 설문 조사이다. 이 설문 자료를 회귀모델로 설정하면 종속변수 Y는 구매 결정의 여부, X_i 는 구매 결정에 미친 항목들이다. 표본회귀방정식은 SPSS로 구한다. 이 결과는 [6]에 있다.

Apriori[3,4]알고리즘을 구현하고 빈발항목집단을 추출하기 위하여 SQL2000 데이터베이스 프로그램과 ASP 방식을 이용하여 구현하였다. 처음 화면에서 각각의 설문 응답자가 보여준 값과 그리고 분석자가 설정하는 threshold값을 입력한다. Apriori[3,4] 알고리즘은 위의 자료를 가지고 입력된 threshold 값에 상응하는 빈발항목집단을 찾아낸다.

3개의 서로 다른 threshold값을 가지고 빈발항목집단을 추출하였다. 그럼 의 설문지에서 “매우 그렇다”에 해당하는 값 5, 그리고 차례로 4, 3을 가지고 빈발항목집단을 추출하였다. 5에 해당하는 결과는 {그래픽}, 4에 해당하는 결과는 {그래픽, 네비게이션 도움}, 그리고 3에 해당하는 결과는 {그래픽, 네비게이션 도움, 색상}으로 결과가 도출되었다.

위의 설문 자료에 해당하는 회귀방정식의 결과는 [6]에 있으면 다음과 같다.

	B	Beta**
constant	.867	
네비게이션 도움	.379	.184
그래픽	.418	.299
색상	2.190E-02	.011
멀티미디어	-8.611E-02	-.040
레이아웃	-7.376E-02	-.033
홈페이지 구조	-.174	-.081

** Beta는 표준화 계수라고 한다. 표준화 계수는 변수가 여러 개 유의할 경우 영향력의 차이를 보는 것이다.

그림4 전자상거래 구매결정에 관한 회귀분석 결과

그림4의 회귀분석결과를 분석하면 구매 결정에 가장 큰 영향을 미치는 순서는 “그래픽”, “네비게이션 도움”, “색상”的 순서임을 알 수 있다. 이는 threshold 값을 큰 값에서 작은 값으로 변화 시키면서 얻은 빈발항목집단 {그래픽}, {그래픽, 네비게이션 도움}, {그래픽, 네비게이션 도움, 색상}과 일치한다. 즉, 이 논문에서 확인하고자 하는 것과 일치한다. 이는 데이터마이닝의 연관관계 추출을 위한 빈발항목집단을 이용하여 회귀방정식 모델을 설정 할 수 있음을 보여준다.

이 논문에서는 다중선형회귀방정식 모델과 빈발항목집단과의 관계를 밝혔다. 그러나 회귀모델에는 선형모델외에도 다양한 모델이 있다. 앞으로 선형모델 이외의 모델의 설정에 도움이 되는 방법의 개발이 요구된다.

참고문헌

- George G. Judge, "Introduction to the theory and practice of econometrics", John Wiley, 1982
- George G. Judge, "The theory and practice of econometrics", John Wiley, 1980
- Ming Syan Chen, Jiawei Han, Philip S. Yu, "Data Mining: An Overview from Database Perspective", IEEE TKDE, vol 8, no 6, 1996, pp 866 - 883.
- Rakesh Agrawal, Tomasz Imielinski, Arun Swami, "Mining Association Rules between sets of items in large databases", Proceedings of the 1993 ACM SIGMOD Conference, May 1993
- 김영인, “이미지에 기반한 패션색채의 데이터베이스 구축 및 실용화 연구”, 산업디자인 기반기술사업 과제 결과 보고서, 1998
- 차영주, 이태경, “인터넷마케팅 전략으로써 쇼핑몰 디자인이 구매결정에 미치는 영향에 관한 연구”, 디자인학연구, vol 14, no 1, 2001, pp 17 ~ 26