

토픽별 인간 속성의 영향력 기반 소셜 관계 지수 산정

권오상*, 박진우*, 이상훈**

*국방대학교 국방정보체계학과

**국방대학교 국방정보체계학과 교수

e-mail:gentle_sea@naver.com

Social Relationship Value Computation based on the Influence of Human Attributes classified by Topics

Oh-Sang Kwon*, Gun-Woo Park*, Sang-Hoon Lee**

Department of Defense Information System, Korea National Defense University

요 약

최근 검색엔진의 효율성을 향상시키고 검색결과에 있어서 사용자들의 요구사항을 충족시키기 위한 연구들이 활발히 수행되고 있으며, 많은 방법론들이 제시되고 있다. 이는 방대한 정보 속에서 사용자의 검색 의도에 맞는 정보를 효과적으로 제공하는 것을 그 목표로 한다. 특히 본 논문에서는 검색하고자 하는 토픽별 사용자의 인적 속성들이 미치는 영향력을 기반으로 사용자간 소셜 관계 지수(SRV : Social Relationship Value)를 산정하는 방법을 제안한다. 소셜 관계 지수란 인간의 내재적인 특성을 수치로 산정한 것으로, 웹 사용자들에게 있어서는 검색 성향의 유사도와 직결된다. 따라서 검색하고자 하는 토픽별 개인 성향의 유사도를 수치로 부여하고 유사성이 높은 사람들의 검색 정보를 이용하면 사용자에게 보다 만족된 검색결과를 제공할 수 있다. 본 연구에서는 구글 디렉터리(Google directory)의 정제된 각 토픽별 하위 범주(category)에 대해 선택 결과가 같은 사람들을 대상으로 인적 속성을 분석하고, 그 영향력을 가중치로 적용해 산정된 소셜 관계 지수와 사용자들의 검색 패턴을 비교 하였다. 그 결과 특정인을 기준으로 소셜 관계 지수가 높은 사람들의 검색 패턴이 매우 유사함을 확인 하였다. 이를 통해 토픽별 개인 간 연결 강도가 강할수록, 즉 유사성이 높은 사용자간에는 검색 패턴 또한 유사함을 검증 할 수 있었다.

1. 서론

오늘날 우리의 일상생활에는 인터넷 서비스를 이용하여 자신의 경험을 공유하거나 필요한 정보를 찾아낸다. 그러나 현재의 웹 환경은 실로 방대한 양의 정보를 포함하고 있어 검색 결과에 대한 사용자의 만족도를 충족시키기란 매우 어렵다. 그래서 최근의 검색엔진들은 사용자 개인에게 집중하는 개인화 검색(Personalized Web Search) 서비스를 제공하고자 노력하고 있다. 즉 검색엔진의 가장 큰 목표는 방대한 정보 속에서 사용자가 원하는 정보를 효율적으로 제공하는 것이다. 본 논문에서는 이러한 어려움을 해결하고자 현재 인터넷에서 활발하게 활용하고 연구가 진행되고 있는 소셜 네트워크(Social Network)와 협업 필터링(CF : Collaborative Filtering) 개념을 검색에 적용하여 토픽별 검색의 효율성과 적합성을 향상시키고자 한다.

본 논문의 구성은 다음과 같다. 2장에서는 관련 연구에 대하여 기술하고, 3장에서는 토픽별 웹 사용자의 인적 속성 영향력을 기반으로 소셜 관계지수를 산정하는 방법을 제시하고, 4장에서는 이를 이용한 웹 검색 실험 및 기존 연구결과와의 비교 결과를 기술한다. 마지막 5장에서는 결론 및 향후 연구를 기술한다.

2. 관련연구

2.1 소셜 네트워크(Social Network)

사회 구성원 간 관계를 맺고 있는 구조는 최근 온라인을 중심으로 하여 소셜 네트워크의 개념으로 소개되고 있다. 이는 하나 이상의 상호의존적인 관계에 의해 구성된 개인 또는 집단으로 사회적 구조체(Social Structure)로 정의할 수 있다. 또한 웹 환경에서 개인 중심의 네트워크로 구성되어 사용자별 프로파일을 탐색하고 새로운 연결 및 정보의 소통을 지원한다. 전 세계적으로 유행하고 있는 대표적인 Social Network Service(SNS)에는 Facebook¹⁾, Orkut²⁾, Myspace³⁾, Twitter⁴⁾ 등이 있다.

소셜 네트워크 분석(Social Network Analysis)은 소셜 네트워크의 형태와 특성을 알고리즘 적으로 연구하는 것으로, 접근방법에는 전체 관계망에서 위치와 그 효과를 측정하는 위치적 접근법, 연결망의 직접적인 관계에 초점을 둔 관계적 접근법으로 분류될 수 있으며, 이는 특히 사회학, 통신공학, 경제학 등에서 폭넓게 연구되고 있다[1][3].

1) www.facebook.com

2) www.orkut.com

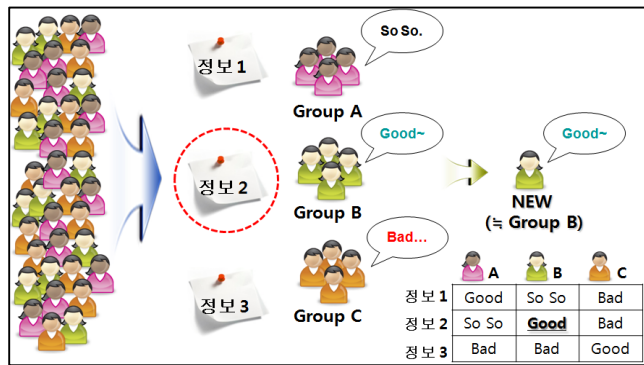
3) www.myspace.com

4) www.twitter.com

2.2 협업 필터링(Collaborative Filtering)

협업 필터링은 서로 비슷한 취향을 가진 사용자간에 아직 접하지 못한 아이템을 교차 추천하거나 분류된 사용자 성향에 알맞은 아이템을 찾아 추천하는 서비스에서 많이 사용되는 방법으로, 보통 많은 사용자들로부터 얻은 기호정보에서 유사도(Similarity)에 따라 그들의 관심사를 자동적으로 예측하게 해주는 개인 추천 서비스 제공 기법이다. 이는 근본적으로 사용자들의 과거 성향이 미래에도 그대로 유지될 것이라는 가정을 전제로 한다. 따라서 단일 사용자가 아닌 다수 사용자들의 과거자료를 바탕으로 공통성(유사도)을 뽑아내고, 이를 분석하여 사용자 선호도에 맞는 새로운 추천을 제공할 수 있는 것이다[5][6].

(그림 1)은 사용자 기반 필터링의 원리를 나타내는 것으로 'Group B'와 유사도가 높은 '신규 사용자(NEW)'에 대해 정보 추천시, 'Group B'가 과거에 만족했던 '정보 2'를 추천 하는 것이 '신규 사용자(NEW)'에게도 역시 만족도가 가장 높다는 것을 의미한다.



(그림 1) 사용자 기반 필터링

3. 토픽별 소셜 관계 지수 산정 방법

개인별 속성이 유사하고 공동의 관심사로 구성되어진 소셜 네트워크를 웹 검색에 적용하면 검색 결과를 상호 추천해줄 수 있다. 왜냐하면 인간의 내재적 특성을 수치화한 소셜 관계 지수가 높다는 것은 웹 검색에 있어 검색 성향의 유사정도가 높다는 것과 같기 때문이다. 그러나 단순히 전체적인 웹 검색 성향이 같다고 해서 모든 토픽에 대해 같은 검색 결과를 만족한다는 것은 아니다. 이는 사용자간 전체적인 소셜 관계 지수가 높아도 토픽이라는 키워드로 세분화 될 때, 토픽별 모든 검색 결과가 동일할 수 없다는 것을 의미한다. 예를 들면, 평소 나와 소셜 관계 지수가 높고 전체적인 웹 검색 성향 또한 유사한 특정 사용자의 경우 토픽 'Game'에 대해서는 강한 소셜 관계를 형성할 수 있지만, 'Art'라는 전혀 다른 분야에 대해서는 소셜 관계가 약하거나 관계 자체가 없을 수도 있다.

따라서 연결 관계를 맺고 있는 웹 사용자들은 모든 토픽에 대해 획일적으로 같은 관심사를 갖고 있지 않으며, 동일한 질의(Query)에 대해 검색 결과의 만족도 또한 서로 다를 수 있다고 가정 한다면, 토픽별로 검색성향에 미치는 인적 속성들의 영향력 정도를 분석하고 이를 기반으로 소셜 관계 지수를 산정 후, 현재와 같은 소셜 네트워크의 검색 환경에 적용 할 때 검색엔진의 효율성과 신뢰성을 보다 향상시킬 수 있을 것이다.

3.1 사용자 프로필

소셜 네트워크상에서 어떠한 토픽에 대해 영향을 미치는 사용자의 프로필을 식별하는 것은 매우 중요하다. 일반적으로 사용자의 프로필에는 성별, 나이, 거주지역 등 밖으로 드러내도 크게 문제가 되지 않는 프로필과 개인의 성격, 재산, 반도덕적 취향 등의 타인에게 드러내고 싶지 않은 프로필로 나뉜다. 연구를 위해서는 사용자의 세부적이고 다양한 프로필이 요구되지만, 개인의 프라이버시 침해라는 문제점 때문에 신중할 필요가 있다.

소셜 관계 지수를 산정하기 위해서는 사용자 프로필 각각의 세부항목을 속성(Attribute)이라 정의하고, 성격이 유사한 범주의 속성들끼리 분류하여 요소(Factor)라고 정의한다. 예를 들어 인적 요소를 유사성(Similarity), 접근성(Access), 친밀성(Intimacy) 등으로 구분한다면, 유사성(Similarity)이라는 요소 안에는 성별, 나이, 출신고 등의 속성들로 구성된다. 이러한 요소와 속성들은 대표적 SNS 사이트인 Facebook, Orkut, Myspace, Twitter 등에서 사용되는 일반적 입력 항목과 웹 환경에서의 검색 행위에서 발생하는 클릭 히스토리 등으로 선정된다.

3.2 토픽(Topic) 선정

사용자의 유사성을 검색 성향에 적용하기 위해 검색 성향을 토픽별로 세분화 한다면, 토픽에 대한 구체적인 정의가 필요하다. 하지만 웹 환경에서 토픽은 큰 범주에서 세부 항목까지 매우 방대하기 때문에 정형화된 형태의 임의적 정의가 어렵다. 그래서 검색 사이트 회사에서는 일반적으로 검색 성능을 향상시키기 위해 모든 검색어를 포함할 수 있는 자회사만의 디렉터리를 분류한 별도의 페이지를 운영한다. 따라서 본 연구에서는 모든 주제를 포함하면서 14개의 커뮤니티로 정제된 구글(Google) 사이트의 디렉터리를 토픽 선정에 이용한다.

3.3 소셜 관계 지수(Social Relationship Value)

사용자들 사이의 기본적인 소셜 관계 지수 산정 방법은 자신이 가지고 있는 인적 속성의 일치 여부로 결정된다. 즉 나와 많은 인적 속성이 일치될수록 강한 소셜 관계를 형성하고 유사도가 높다는 것이다.

<표 1> 속성에 대한 사용자간 Mapping Matrix

Factor	Similarity					Access				...	
Attribute	s1	s2	...	sm	Sum (Sm)	a1	a2	...	an	Sum (An)	...
User _{-me}	1	1	...	1	m×1	1	1	...	1	n×1	...
User 1	1	0	...	1	$\sum_{i=1}^m si$	0	1	...	0	$\sum_{j=1}^n aj$...
User 2	0	1	...	0	$\sum_{i=1}^m si$	1	1	...	0	$\sum_{j=1}^n aj$...
i	i	i	i	i	i	i	i	i	i	i	...

<표 1>은 특정인 즉 나(User_{-me})를 기준으로 다른 사용자들이 가진 속성의 일치여부에 따라 1과 0의 값으로 표현한 Mapping Matrix이다. 그리고 웹 사용자 간의 소셜 관계 지수 산정 알고리즘은 수식(1)과 같다.

$$SRV_{me_user} = \alpha \cdot \frac{\sum_{i=1}^m si}{SI} + \beta \cdot \frac{\sum_{j=1}^n aj}{AI} + \dots \quad (1)$$

- α, β, \dots : 가중치(Balance Factors) / $\alpha + \beta + \dots = 1$
- si, aj, \dots : 웹사용자간 일치하는 속성들
- SI, AI, \dots : 각요소별 속성들의 개수

수식(1)에서는 각 요소에 대해 가중치를 주어 각 요소별 상대적 영향력을 달리하였으며, 각 속성들의 일치여부에 따라 부여한 1과 0의 총합을 각 요소별 속성들의 개수로 나누어 평균화함으로써 요소별 속성의 개수 차이를 해결하였다. 그러나 유사성, 접근성, 친밀성 등과 같이 ‘요소’의 의미가 추상적이어서 가중치를 각 요소별로 주는 것은 적절치 못하다. 또한 각 요소별 포함된 속성들의 상대적 가중치는 고려하지 못한 문제가 발생한다. 따라서 소셜 관계 지수 산정에 중요한 역할을 하면서도 구체적 사실적 의미인 ‘속성’에 가중치를 주는 것이 합리적이다. 즉 성별, 나이 등과 같이 각 속성별 상대적 가중치를 준다면 추상적 요소의 범주에 구애받지 않고 소셜 관계 지수 산정 알고리즘을 보다 간단한 수식(2)로 표현이 가능하다.

$$SRV_{me_user} = \sum_{i=1}^m si(\alpha + \beta + \dots) + \sum_{j=1}^n aj(\gamma + \delta + \dots) + \dots \quad (2)$$

- $\alpha, \beta, \gamma, \delta, \dots$: 가중치(Balance Factors) / $\alpha + \beta + \dots + \gamma + \delta + \dots = 1$
- si, aj, \dots : 웹사용자간 일치하는 속성들

3.4 토픽별 속성의 영향력

수식(2)에서 사용된 속성별 가중치($\alpha, \beta, \gamma, \dots$)는 소셜 관계 지수를 산정함에 있어 사용자 간의 유사성향 정도를 결정하는 각 속성별 영향력과 같다. 예를 들어 친구 맺기의 경우 성별, 나이, 관심분야 등 어느 속성이 더 중요하고 덜 중요한가에 대해 각 속성별 영향력이 가중치로 작용하는 것이다. 본 연구에서는 사용자들이 특정 토픽에 대한 검색에 있어서 각 속성이 미치는 영향력을 구하여 소셜 관계 지수를 산정하는 가중치로 적용한다[2][4].

<표 2> 토픽별 인적속성의 영향력(가중치) Matrix

Attribute User- <i>TOPIC</i>	Attribute 1	Attribute 2	Attribute 3	...	Sum (An)
User- <i>TOPIC1</i>	α_1	β_1	γ_1	...	1
User- <i>TOPIC2</i>	α_2	β_2	γ_2	...	1
⋮	⋮	⋮	⋮	⋮	⋮
User- <i>TOTAL</i>	α	β	γ	...	1

<표 2>는 토픽별로 세분화된 인적속성의 영향력들과 전체 영향력의 관계를 나타내는 Matrix로서 ‘Attribute 1’의 전체 영향력(α)은 토픽별로 세분화된 영향력들($\alpha_1, \alpha_2, \dots$)의 평균과 같다. 따라서 전체 영향력의 크기가 ‘ $\alpha < \beta < \gamma$ ’의 관계가 성립하더라도, 세부 토픽별로 볼 때 ‘ $\alpha_1 > \beta_1 > \gamma_1$ ’와 같이 각 토픽의 특성에 따라 상반된 영향력(가중치)의 결과가 나올 수도 있다.

4. 실험 및 방법

4.1 데이터 셋(Data Set)

본 연구에서는 토픽별 인적속성의 영향력을 기반으로 소셜 관계 지수를 산정하고 이것이 검색 성향에 있어서도 직접적인 관련성이 있는지 알아보려고 한다. 그래서 실험을 위한 토픽의 선정은 대표적 검색 사이트인 ‘구글’의 디렉터리를, 인적속성은 ‘SNS’의 사용자 프로필 항목을 정제하여 사용하였다. 그리고 사용자 프로필의 항목을 고려 다양한 직업군, 학력, 연령층 등의 대상자 700명을 대상으로 개인의 프로필과 토픽별 세부 카테고리들 선택하도록 하였다. (그림 2)는 실험에 사용된 설문지로 개인의 프로필 작성과 토픽별 관심 있는 세부 카테고리 항목들을 선택하도록 구성되어 있다.

성 별	관 계	연령대	출신지역	거주지역	학 력	직업군	온라인 커뮤니티 사이트상 주
남() 여()	기 혼() 연애중() 솔 로()	10대()	강원도()	강원도()	고졸이하()	고위관리() 권문직() 사무()	친구 사귀기 취미활동 개인 인맥 만남/데이트
		20대()	경기도()	경기도()	대학재학()	기술() 판매() 서비스()	
		30대()	경상도()	경상도()	대학졸업()	농/임/어업() 가사()	
		40대()	전라도()	전라도()	대학원재학()	군인() 단순노무() 학업()	
		50이상()	제주도()	제주도()	대학원졸업()	무직() 기타()	
주제별 카테고리							
카 테 고 리							
예 술	애니메이션, 유희품, 건축, 미술의 역사, 신체예술, 고전연구, 만화책, 의상, 공연물, 댄스, 디지털, 인터랙티브, 그래픽 디자인, 인문, 삽화, 문학, 영화, 음악, 신화/전설, 부속, 온라인공연, 사진, 라디오, 수사학, 텔레비전, 극장, 타이포그래피, 비디오, 시각예술, 작각 자						
비즈니스	회계, 비즈니스와 사회, 협동조합, 고객 서비스, 전자상거래, 교육/연수, 취업, 인적자원, 경						

(그림 2) 개인 프로필 및 토픽별 세부 카테고리

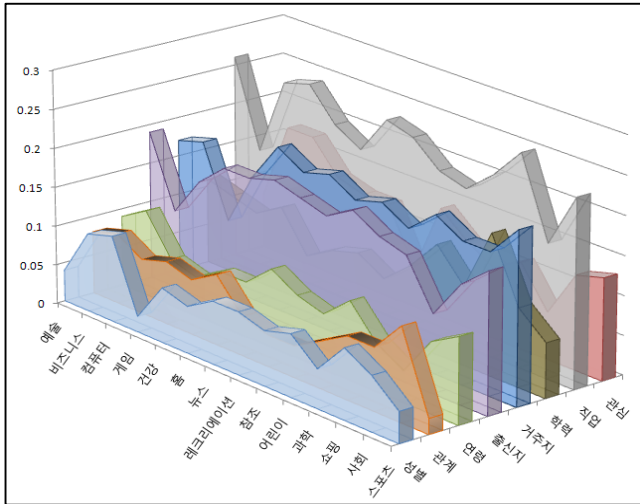
4.2 토픽별 인적속성의 영향력(가중치) 산정

설문을 통해 생성된 700명의 자료에서 토픽별 사람들이 가장 많이 선택한 세부 카테고리 상위 30% 항목을 선별 후 선별된 항목을 기준으로 동일한 선택을 한 사람들을 검색성향이 유사한 것으로 보고, 선택 비율에 따라 검색성향 유사도의 수준을 ‘상·중·하’로 인원을 재분류한다. 이렇게 재분류된 그룹을 대상으로 각 개인 속성별 일치율을 이용해 상대적 영향력을 산출한다. 이때 방법은 설문 결과 각 속성별 세부사항에서 선택된 전체 비율을 ‘영향력 0’로 기준하고, 그 기준에서 벗어난(치우친) 정도, 즉 표준편차를 이용하여 값을 구한다. 그 결과 <표 3>과 같이 토픽별 인적속성의 상대적 영향력을 산정할 수 있다

<표 3> 토픽별 인적속성 영향력 산정 결과

Attribute User- <i>TOPIC</i>	성 별	관 계	연 령	출신 지역	거주 지역	학 력	직 업	관 심	계
예 술	0.04	0.08	0.09	0.19	0.17	0.11	0.26	0.07	1.01
비즈니스	0.1	0.09	0.11	0.1	0.18	0.1	0.15	0.17	1
컴퓨터	0.11	0.07	0.06	0.15	0.09	0.09	0.25	0.17	0.99
게 임	0.02	0.08	0.05	0.18	0.16	0.11	0.26	0.13	0.99
건 장	0.07	0.07	0.07	0.18	0.21	0.06	0.22	0.12	1
홈	0.06	0.09	0.07	0.19	0.19	0.08	0.2	0.12	1
뉴스	0.08	0.05	0.1	0.18	0.2	0.09	0.25	0.05	1
여가활동	0.08	0.04	0.08	0.17	0.18	0.07	0.24	0.14	1
참 조	0.07	0.04	0.07	0.19	0.19	0.1	0.21	0.11	0.98
지 역	0.08	0.06	0.1	0.17	0.17	0.14	0.2	0.09	1.01
과 학	0.05	0.08	0.06	0.16	0.2	0.1	0.23	0.11	0.99
쇼 핑	0.09	0.08	0.04	0.1	0.18	0.18	0.27	0.06	1
사 회	0.07	0.12	0.09	0.15	0.18	0.1	0.17	0.12	1
스포츠	0.04	0.02	0.11	0.18	0.22	0.07	0.24	0.13	1.01
User- <i>TOTAL</i>	0.07	0.07	0.08	0.16	0.18	0.1	0.23	0.11	1.01

(그림 3)은 토픽에 대해 영향을 미치는 인적속성의 가중치 변화를 도식화한 그래프로 특정 토픽에 따라 인적속성의 영향력이 변화한다는 것을 나타낸다.



(그림 3) 토픽별 인적 속성 가중치 변화

4.3 토픽별 소셜 관계 지수 산정

<표 3>에서 구한 영향력을 수식(2)의 가중치 값으로 넣고, 속성에 대한 사용자간 Mapping Matrix인 <표 1>에 적용하면, 사용자별 전체 소셜 관계 지수인 <표 4>를 구할 수 있으며, 또한 토픽별 소셜 관계 지수인 <표 5>로도 재구성하여 표현할 수 있다.

<표 4> 사용자별 전체 SRV

Attribute User	성별	관계	...	관심	계
User-me	1	1	...	1	1
User 1	0.07	0	...	0.11	0.43
User 2	0	0.07	...	0.11	0.56
User 3	0.07	0.07	...	0	0.51
⋮	⋮	⋮	...	⋮	⋮

<표 5> 토픽별 SRV

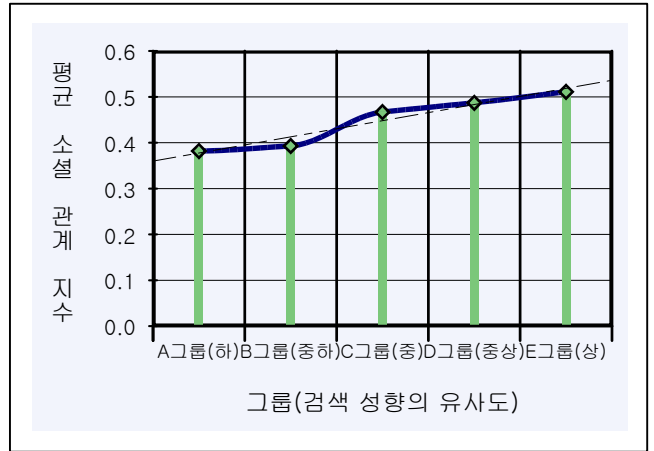
Topic User	예술	비즈니스	...	스포츠	평균
User-me	1	1	...	1	1
User 1	0.49	0.38	...	0.51	0.43
User 2	0.48	0.59	...	0.49	0.56
User 3	0.52	0.51	...	0.47	0.51
⋮	⋮	⋮	...	⋮	⋮

<표 5>에서 보듯이 특정인인 '나'에 대해 'User 2'의 전체적인 소셜 관계 지수가 가장 높지만, 토픽별로 볼 때 '예술'에서는 'User 3', '스포츠'에서는 'User 1'의 소셜 관계 지수가 가장 높다. 따라서 전체적인 소셜 관계 지수가 높은 사용자라고 해서 모든 토픽에 대해서도 소셜 관계 지수가 항상 높은 것은 아니라는 것을 알 수 있다.

4.4 토픽별 소셜 관계 지수-검색 성향 간 상관관계 분석

소셜 관계 지수 정도가 웹 검색 성향의 유사도와 연관성이 있는지를 검증하기 위해 토픽별 '나'와 같은 카테고리를 선택한 사용자들을 재분류하여 검색 성향 유사 정도에 따라 '하·중하·중·중상·상'의 5개 그룹(A~E)으로 나누었다. 그리고 '나'에 대한 각 그룹별 구성원들의 소셜 관계 지수 평균값을 비교하였다. 그 결과 검색 성향 유사도가 가장 높은 'E그룹'의 평균 소셜 관계 지수 값은 0.511로 가장 높고, 검색 성향 유사도가 가장 낮은 'A그룹'은 0.382로 가장 낮은 결과 값을 얻었다. 이를 도식화한 (그림 4)는 검색성향의 유사도에 따라 그룹별 평균

소셜 관계 지수 값의 변화를 나타내는 것으로 평균검색성향의 유사도가 높은 그룹일수록 평균 소셜 관계 지수가 높아지는 경향을 확인할 수 있다.



(그림 4) 검색성향 유사도에 따른 평균 소셜 관계 지수

5. 결론 및 향후 연구

방대한 양의 정보를 포함하고 있는 웹 환경에서 검색 엔진의 효율성과 신뢰성을 향상시키기 위한 연구들이 활발히 진행되어 왔다.

본 논문에서는 웹 사용자 간의 소셜 관계 지수를 토픽별 인적 속성의 영향력을 이용하여 산정 후, 검색 성향과의 상관관계를 분석하였고, 그 결과 현재와 같은 소셜 네트워크 검색 환경에 적용 시 검색엔진의 효율성과 신뢰성을 보다 향상시킬 수 있다는 가능성을 확인하였다.

향후 과제로 보다 균등하고 다양한 개인 프로파일 수집 및 그들의 검색 성향 파악 등 더 많은 인원을 대상으로 실험하여 연구의 효과성을 진단할 것이다. 또한 토픽과 속성간의 일반적인 상관관계, 즉 연관성 여부 및 연관 수준 값을 측정하여 객관적인 토픽-속성 간 연관성 Matrix를 생성한다면 보다 정확하고 신뢰할 수 있는 토픽별 소셜 관계 지수 산정이 가능할 것이다.

참고문헌

- [1] YongHak Kim, "Social Network Analysis", ParkYoung Company, 2003
- [2] John Holt, "NEC social and organizational factors", HVR Consulting Services Ltd, 2003
- [3] JiSu Kim, "Human Network", Policy of Information and Communication, a volume 16, an issue 16, 2004
- [4] AnthonyH.Dekker, "Revisiting SCUDHunt and the Human Dimension of NCW:Some Thoughts", Defence Systems Analysis Division, DSTO, 2006
- [5] Wei Chu, SeungTaek Park, "Personalized Recommendation on Dynamic Content Using Predictive Bilinear Models", WWW 2009 MADRID
- [6] WenYen Chen, JonChyuan Chu, Junyi Luan, "Collaborative Filtering for Orkut Communities:Discovery of User Latent Behavior", WWW 2009 MADRID