

구글 검색엔진을 활용한 키워드 검색결과 수 관리 시스템 설계 및 구현⁺

이주연 · 이중화 · 박유현*

A design and implementation of the management system for number of keyword searching results using Google searching engine

Ju-Yeon Lee · Jung-Hwa Lee · Yoo-Hyun Park*

Department of Computer Software Engineering, Dongeui University, Busan 47227, Korea

요 약

인터넷 상에 많은 정보들이 발생하면서 검색 엔진은 사용자에게 필요한 흩어진 정보를 모아주는 중요한 역할을 하고 있다. 일부 검색 엔진에서는 검색어가 포함된 검색 결과 페이지뿐만 아니라 검색 결과 수도 함께 제공하고 있다. 구글 검색엔진에서 제공하는 검색 결과 수는 인터넷에서 해당 검색어에 대한 전체적인 추세를 파악하는데 활용될 수 있다. 본 논문에서는 구글 검색엔진에서 제공하는 검색결과 수를 효과적으로 관리할 수 있는 구글 검색엔진을 활용한 키워드 검색결과 수 관리 시스템을 설계하고 구현하고자 한다. 제안하는 시스템은 웹으로 작동하며 검색 에이전트, 저장 노드, 검색 노드로 구성되어 키워드 및 검색 결과 수를 관리하고 검색을 수행한다. 최종 검색 결과로는 검색 키워드, 검색 결과 수, 검색 결과 수를 활용하여 두 키워드의 거리를 계산하는 NGD(Normalized Google Distance)가 제공된다.

ABSTRACT

With lots of information occurring on the Internet, the search engine plays a role in gathering the scattered information on the Internet. Some search engines show not only search result pages including search keyword but also search result numbers of the keyword. The number of keyword searching result provided by the Google search engine can be utilized to identify overall trends for this search word on the internet. This paper is aimed designing and realizing the system which can efficiently manage the number of searching result provided by Google search engine. This paper proposed system operates by Web, and consist of search agent, storage node, and search node, manage keyword and search result, numbers, and executing search. The proposed system make the results such as search keywords, the number of searching, NGD(Normalized Google Distance) that is the distance between two keywords in Google area.

키워드 : 검색 엔진, 병렬 시스템, 키워드 검색, 빅데이터, 데이터과학

Key word : Search Engine, Parallel system, keyword search, Big Data, Data Science

Received 16 February 2016, Revised 22 February 2016, Accepted 16 March 2016

+ 본 논문은 동의대학교 석사학위 논문 "확장 가능한 병렬 키워드 검색 시스템 설계"의 내용을 요약 정리하여 작성하였음

* Corresponding Author Yoo-Hyun Park(E-mail:yhpark@deu.ac.kr, Tel:+82-51-890-1737)

Department of Computer Software Engineering, Dong-Eui University, Busan 47227, Korea

Open Access <http://dx.doi.org/10.6109/jkice.2016.20.5.880>

print ISSN: 2234-4772 online ISSN: 2288-4165

©This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License(<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.
Copyright © The Korea Institute of Information and Communication Engineering.

I. 서론

인터넷이 발전함에 따라 인터넷 상에는 많은 정보들이 만들어지고 저장되어 있어 사용자가 이를 개별적으로 관리하는 데는 문제가 많다. 이러한 이유로 인터넷 상의 수많은 정보들 중에서 사용자들에게 필요한 정보만을 선별해 주는 검색엔진은 매우 유용한 도구가 되어 왔다. 현재 검색엔진의 점유율은 국내에서는 네이버, 다음이, 중국에서는 바이두가 가장 높으며, 전 세계적으로는 구글(google)이 압도적으로 많이 사용되고 있다.

한편, 최근 많은 분야에서 연구하고 있는 빅데이터는 3V(Volume, Velocity, Variety)의 특성을 가지는 데이터를 의미하며 기존의 처리방법으로는 이러한 데이터를 분석하는데 많은 문제가 있어, 이에 맞는 새로운 분석 방법들이 연구되고 있다. 인터넷 역시 빅데이터의 대상으로 볼 수 있으며, 인터넷 상의 빅데이터를 분석하는 가장 쉬운 접근법 중 하나가 구글의 검색 결과 수를 활용하여 전체적인 추세를 분석하는 것이다.

구글 검색 엔진은 어떤 페이지가 다른 페이지에 얼마나 참조되었는가를 판단하는 페이지랭크(PageRank) 기술[1]을 사용하여 검색 결과를 제공 하는데, 현재 서비스를 제공하는 대부분의 검색엔진과는 달리 검색 키워드가 포함된 검색 결과 페이지뿐만 아니라 ‘검색 결과의 수’도 함께 보여주고 있다. 이러한 검색 결과의 수는 그 활용방법에 따라 다양한 분야에서 활용될 수 있다. 특히, 미국의 독감환자 통계[2]는 대표적인 검색결과 수를 활용한 예이며, 그밖에도 미국 대선 예측, 한국 대선 예측, 뉴욕의 연방준비은행의 미국의 주택지표, 달러, 위안, 환율, 독일 실업률과 경제지표 예측에도 구글 검색 결과 수를 활용한 예가 있다.

만일 대선 후보 예측과 같이 특정 시점에서 특정 키워드를 검색한 결과의 수가 필요한 경우에는 비교적 간단한 검색만으로도 이러한 정보를 제공할 수 있겠지만, 보다 많은 키워드를 검색하거나 키워드 쌍을 조합하여 검색해야 하는 경우에는 이를 위한 전용 시스템을 사용하는 것이 보다 효율적이다.

본 논문에서는 복수개의 키워드를 동시에 검색하여 그 검색결과를 자동으로 추출하고 이를 관리할 수 있는 할 수 있는 구글 검색엔진을 활용한 키워드 검색결과 수 관리 시스템을 설계하고 구현하고자 한다. 본 논문

에서 제안하는 시스템은 크게 검색 에이전트, 저장 노드, 검색 노드로 구성되어 동작하며, 최종 검색 결과로는 검색 키워드, 검색 결과 수, 검색 결과 수를 활용하여 두 키워드의 거리를 계산하는 NGD(Normalized Google Distance)[3]가 제공된다.

II. 관련연구

2.1. 빅데이터 빈도수를 이용한 추세 분석 사례

구글에서는 특정 단어에 대한 시간적인 변화를 구글 트렌드[4]라는 서비스로 제공하고 있다. 미국의 독감환자 통계[2]는 가장 대표적인 검색결과 수를 활용한 예이고, 이는 보건 의료 분야에서 빅데이터(Big Data)의 활용 가능성을 보여주는 중요한 사례이다. 구글은 2008년 11월부터 독감 트렌드 서비스를 시작하였는데, 전 세계 각지에서 ‘독감증세’, ‘독감치료’ 등 독감과 관련된 검색어의 입력 빈도를 지역별로 파악해 독감 유행 수준을 ‘매우 낮음’부터 ‘매우 높음’까지 5개 등급으로 구분해 표시한다[5]. 특정 지역에서 발열이나 기침 등 독감 관련 검색이 늘어나면 검색어와 관련된 IP주소를 지도에 추가해 해당 지역의 독감 유행 수준 등급이 거의 실시간으로 표시된다[5]. 구글은 미국 CDC의 관련 보고서보다 1주에서 2주 정도 더 빨리 독감 바이러스의 활성을 정확히 예측하는 실시간 감시 시스템으로 변환시켜 주는 컴퓨터 모델을 제시했다[6]. 이와 유사하게 트위터를 이용하여 독감을 예측하는 방법도 연구되었다[6].

또한, 구글은 3000만여 권 이상의 책을 OCR(optical character reader)을 통해 문자화 작업을 하여 구글 엔그램 뷰어[7] 서비스를 제공하고 있다. 이 서비스를 통해 디지털 작업을 한 책들 내에 있는 8,000억 개 정도 되는 단어의 사용 빈도 추이를 그래프로 보여주고 있다.

2.2. NGD(Normalized Google Distance)

구글의 독감 트렌드 서비스는 2013년 예측 내용과 다르게 예측되어 정확성이 논란이 있고, [8]의 연구에서는 초콜릿 소비량과 노벨 수상자 숫자 간에 상관관계가 존재하고 있다고 주장하였다. 하지만 이들은 빅데이터는 현상의 모습은 보여줄 수 있지만, 이에 대한 인과관계를 보여주는 데는 한계가 있음을 간과한 예로 볼 수 있다. 이러한 문제들을 해결하여 보다 정확한 추세를

알기 위해서는 관련이 있는 여러 단어들을 함께 검색하는 방법을 고려할 수 있다.

관련이 있는 단어들의 쌍을 찾기 위해서는 그들 간의 관계를 수치화하여 이 수치 값을 비교하는 방법이 유용하다.

두 단어 간의 연관도를 측정하는 방법 중에 정규화된 구글 거리를 나타내는 NGD(Normalized Google Distance)[3]는 두 개로 구성된 객체(단어 또는 구)의 쌍이 가진 정보 거리를 측정하는 개념이며, <식 1>은 NGD 값을 계산하는 식이다.

$$NGD(x,y) = \frac{\max\{\log f(x), \log f(y)\} - \log f(x,y)}{\log M - \{\min\log f(x), \log f(y)\}} \quad (1)$$

(1)에서 $f(x)$ 는 2개의 객체 중 첫 번째 객체를 검색한 검색 결과 수를, $f(y)$ 는 두 번째 객체를 검색한 검색 결과 수를 의미하며, $f(x,y)$ 는 두 객체를 공백으로 구분하여 검색 엔진에서 검색한 검색 결과 수를, M 은 구글의 전체 검색 결과 수를 의미한다.

III. 제안된 키워드 검색결과 수 관리 시스템 설계

본 논문에서는 구글 검색엔진을 활용하여 개별 키워드의 검색결과 수뿐만 아니라 복수개의 키워드를 동시에 검색한 결과 수까지도 관리할 수 있는 구글 검색엔진을 활용한 키워드 검색결과 수 관리 시스템을 제안한다.

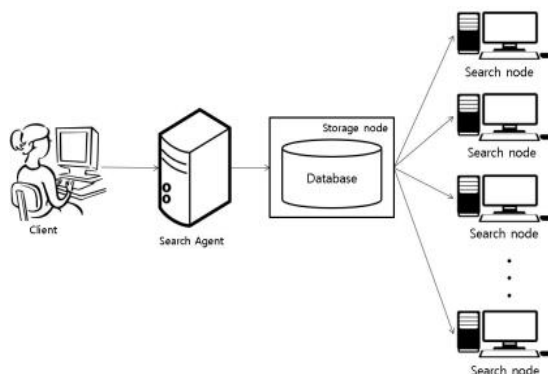


Fig. 1 System Architecture

그림 1은 본 논문에서 제안하고 있는 키워드 검색결과 수 관리 시스템의 구조도이다. 그림 1과 같이 제안하는 시스템은 검색 에이전트, 저장 노드, 검색 노드로 구성되어 있다.

검색 에이전트는 키워드 입력 및 검색 노드의 실행, 결과 화면 출력 등의 역할을 한다. 저장 노드는 검색 키워드 및 검색 결과 수 등의 검색 정보를 저장하는 역할을 수행하는 데이터베이스를 의미하며, 검색 노드는 저장 노드에 저장된 검색 키워드를 실제 구글에 검색을 요청하는 노드들이다.

3.1. 검색 에이전트

검색 에이전트는 제안하는 시스템에서 사용자가 원하는 키워드를 저장 노드에 입력하고, 검색 노드의 검색 수행, NGD 값 계산 최종 검색 결과를 사용자에게 제공하는 역할을 한다.

먼저 사용자로부터 검색 키워드를 입력 받는 방법은, 사용자로부터 직접 키워드를 입력 받거나 TXT 파일의 형태로 입력받을 수 있다. 입력받은 키워드를 기반으로 검색 에이전트는 키워드를 공백을 사이에 두고 2개의 쌍으로 조합하여 생성한다. 키워드 조합은 저장 노드에 저장되며, 이들 키워드 조합의 검색 결과는 다양한 형태로 사용자에게 제공된다.

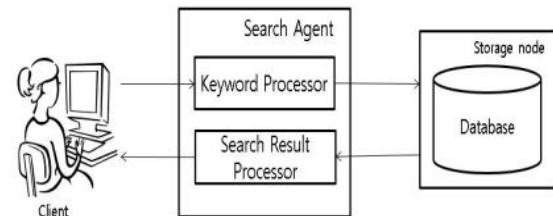


Fig. 2 Search Agent Architecture

검색 에이전트는 그림 2와 같이 키워드 처리기와 검색 결과 처리기로 구성된다.

키워드 처리기는 검색할 키워드 리스트를 입력받고 이를 통해 검색할 키워드 쌍을 생성하는 역할을 수행한다. 즉, 사용자 환경을 통해 직접 키워드들을 입력 받거나 파일의 형태로 입력된 키워드들은 그림 3과 같이 키워드 쌍이 생성되어 저장노드의 데이터베이스에 저장된다.

검색 결과 처리기는 검색 노드에서 검색을 수행한 후

에 저장노드에 입력된 검색 결과 값을 활용하여 사용자에게 최종적으로 다양한 형태의 결과를 보여주는 역할을 수행한다. 본 논문에서는 두 키워드의 정보거리인 NGD값을 제공한다. 본 논문에서는 NGD를 구하는 식 (1)에서 M의 값을 [3]에서 최대범위로 보았던 $9 \cdot 10^9$ 로 가정하여 계산하였다.

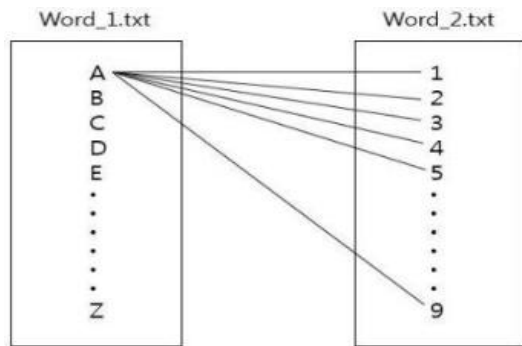


Fig. 3 Generate Keyword pairs with files

3.2. 저장 노드

저장 노드는 검색 키워드 쌍과 이에 관한 여러 정보를 관리한다. 검색 키워드 쌍에 관한 정보에는 키워드 쌍의 내용은 물론 각 키워드들의 검색 결과 수, 키워드 쌍을 함께 검색한 결과 수가 포함된다. 아울러 검색노드에 아직 검색이 완료되지 않은 검색 키워드 쌍을 찾을 수 있도록 하는데 필요한 여러 정보도 관리한다. 저장 노드에서는 이러한 정보를 위해 테이블 형태의 데이터베이스를 유지하며, 이에 대한 구체적인 구조는 표 1과 같다.

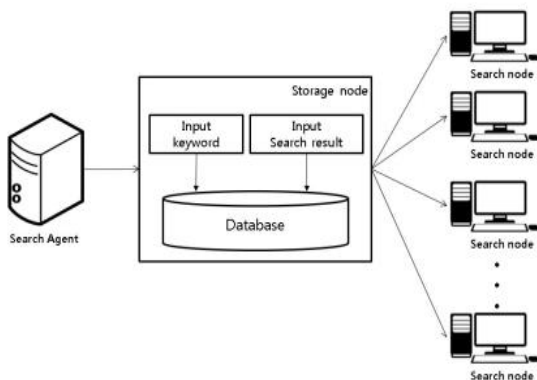


Fig. 4 Design of Storage node

그림 4는 저장노드의 구조를 나타낸 그림이다. 그림과 같이, 저장 노드는 키워드 입력기와 검색 결과 입력기로 구성되어 있다.

키워드 입력기는 검색 에이전트가 생성한 키워드 쌍을 전달 받아 데이터베이스에 초기 값을 저장하는 역할을 한다. 새로운 키워드 쌍 검색을 위해서는 먼저 데이터베이스를 초기화 하여 이전에 사용된 키워드 쌍 정보를 삭제한 후 새로운 키워드 쌍을 저장해야만 하며, 키워드 쌍을 제외한 모든 속성 값은 초기화 된 값이 저장된다.

저장노드에서 사용하고 있는 테이블 형태의 데이터베이스는 표 1과 같이 10개의 속성을 관리한다. SearchWord는 검색 에이전트로부터 전달받은 키워드 쌍을 의미하며, 그 각각의 키워드들은 Keyword1과 Keyword2에 저장된다. 또한 Result는 SearchWord(키워드 쌍)으로 구글 검색한 결과 수이며, Result1과 Result2는 각각 Keyword1과 Keyword2의 구글 검색 결과 수를 저장한다. NGD는 Keyword1과 Keyword2의 정보거리를 <식 1>에 의해 계산한 결과가 저장된다. Who, GetTime, InputTime 등은 여러 검색 노드들 중에서 어떤 노드가 해당 키워드 쌍을 검색할 것인지 등을 결정하는데 사용되는 속성 값이다. 이에 대해서는 3.5에서 자세히 설명한다.

Table. 1 Attribute of Storage node

Attribute	Explanation
SearchWord	Keyword pair
get_time	Time of taking search word to the search node
result	Search result number of SearchWord
who	Computers name
input_time	Input time of result number
keyword_1	First Keyword of Searchword
keyword_2	Second Keyword of Searchword
result_1	Search result number of keyword_1
result_2	Search result number of keyword_2
NGD	NGD number

저장 노드의 검색 결과 입력기는 검색 노드에서 검색을 통해 추출한 검색 결과를 저장 노드에 저장하는 역할을 한다. 즉, result(키워드 쌍의 검색결과 수), Result1(Keyword1의 검색결과 수), Result2(Keyword2의 검색

결과 수)의 실제 값이 이때 저장된다.

3.3. 검색 노드

저장 노드에 입력된 검색 키워드 쌍에 대한 실질적인 대용량 키워드 검색은 시스템의 검색 노드들에서 수행한다. 검색 노드는 서로 다른 인터넷 주소를 사용하는 다수의 노드들로 구성되며, 검색 에이전트를 통해 저장 노드에 저장된 키워드 쌍에 대한 검색을 수행하는 검색 명령이 내려진다.

검색 에이전트로부터 검색 명령을 받은 검색 노드들은 SQL 쿼리에 지정된 LIMIT 수만큼의 키워드 쌍에 대한 검색을 수행한다. 본 논문에서는 한 번의 검색에 50개의 키워드에 대한 검색을 수행한다. 검색 노드는 검색을 위한 키워드를 불러올 때 저장 노드에서 관리하는 기본적인 속성에 공백이 있을 경우 검색을 중지하고 임의의 시간을 대기한 후 다시 검색을 수행한다. 검색 노드는 데이터베이스에 표 2와 같이 질의하여 검색할 대상 키워드 쌍의 정보를 가져온다.

즉, 어떤 검색노드도 해당 키워드 쌍을 검색 중이지 않은 키워드 쌍을 가져가게 되며, 임의의 노드가 해당 키워드 쌍을 검색하도록 결정된 후에는 해당 키워드 쌍의 Who에는 검색노드의 노드의 이름을, GetTime에는 키워드 쌍을 가져간 저장 노드에서의 시간을 저장한다.

Table. 2 Query example for deciding the searching keyword pair

```

SELECT SearchWord
FROM InfoTable
WHERE Who != '' and GetTime != ''
ORDER BY Keyword1 LIMIT 50;
    
```

검색이 끝나면 해당 키워드 쌍의 검색결과 수인 Result, 각각의 키워드에 대한 검색결과 수인 Result1, Result2의 값을 저장한다. 또한, 검색이 종료되었음을 나타낼 수 있도록 InputTime에 검색노드에서 저장노드에 입력한 Result 등의 값이 저장된 시간을 저장한다.

이러한 검색 과정이 끝나면 새로운 검색 키워드 쌍을 가져가기 위해 표 2의 과정부터 다시 수행하며, 해당 질의의 결과가 없을 때까지 반복한다. 그림 5는 클라이언트가 검색 에이전트를 이용하여 저장노드, 구글 검색

엔진을 활용하여 키워드 쌍을 검색하는 검색 노드의 검색 과정을 나타낸 그림이다.

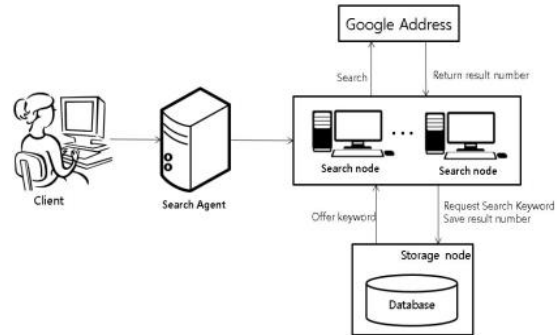


Fig. 5 Search Process

검색 에이전트에서 검색 명령을 내리면 검색 노드에서는 실제 검색을 위한 검색 프로그램이 동작하면서 검색이 이루어지며 본 논문에서는 PHP 스크립트로 구현하였다. PHP 스크립트에서는 검색 키워드를 PHP의 urlencode()함수를 이용하여 인코딩 한 후 검색 엔진에서 검색 키워드가 들어가는 특정 위치에 입력하여 검색한다. 검색을 수행한 후 검색 결과 화면에서 검색 결과 수를 의미하는 '검색결과 약' 과 '개'사이의 숫자를 파싱하여 검색 결과 수로 저장 노드의 Result속성에 저장된다. 표 3은 PHP 스크립트에서 검색 키워드를 인코딩한 후 인코딩된 검색 키워드가 들어가는 자리를 나타낸 표이다.

Table. 3 Encoding keyword position

```

$url = "http://www.google.com/search?q=
    search&btnG=Search&meta=&ie=utf-8&oe=utf-8";
    
```

IV. 제안된 키워드 검색결과 수 관리 시스템 구현

본 논문에서 설계한 시스템은 리눅스 환경에서 구현되었다. 제안하는 시스템에 사용한 운영체제는 UBUNTU 14.04로 검색 에이전트, 검색노드, 저장노드에 모두 사용되었다. 저장노드의 검색 키워드 및 검색 결과 수를 관리하는 데이터베이스로 MYSQL-SERVER

5.5, 검색을 수행하는 프로그램은 PHP5, JAVASCRIPT, HTML5를 사용하여 구현하였다. 본 논문에서는 검색에 8대의 검색 노드를 사용하였다. 그림 6은 웹을 통해 접근하는 시스템 초기화면이다.

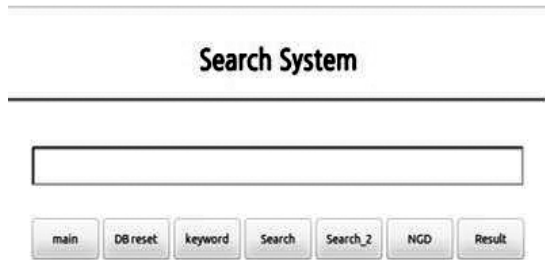


Fig. 6 System Screen in the beginning

시스템의 동작은 초기화면의 버튼이 배치되어 있는 순서대로 진행된다. 검색 키워드를 입력하기 전에 저장 노드의 데이터베이스를 초기화 하고(DB 초기화 버튼), 파일을 이용하거나 사용자가 직접 키워드를 입력하여 검색 키워드를 저장노드에 저장한다(키워드 입력 버튼). 저장 노드에 검색 키워드가 입력되면 검색 버튼을 통해 검색 노드에서 실제 검색이 이루어진다. 저장 노드에 저장된 검색 결과 값으로 NGD 값을 계산하고 결과 보기 버튼을 통해 저장 노드의 키워드, 검색 결과 수, 검색 결과 수를 활용하여 계산한 NGD 값이 최종 검색 결과로 클라이언트에게 제공된다. 최종 검색 결과는 검색 결과 순으로 출력하거나 NGD 순으로 출력할 수 있다.

SQL	FACEBOOK	2015-11-22 20:17:39	72000000	915-8	2015-11-22 20:18:26	SQL	FACEBOOK	2350
00000	2147483647	0.27150131226094						
SQL	ONCS	2015-11-22 20:17:39	2580000	915-8	2015-11-22 20:18:27	SQL	ONCS	2350
00000	4740000	-1.537307305871						
SQL	H4000P	2015-11-22 20:17:39	5830000	915-8	2015-11-22 20:18:28	SQL	H4000P	2350
00000	22300000	-8.41897567948888						
SQL	HPD	2015-11-22 20:17:39	2830000	915-8	2015-11-22 20:18:30	SQL	HPD	2350
00000	38300000	-1.2833436524100						
SQL	INTEL	2015-11-22 20:17:39	9670000	915-8	2015-11-22 20:18:30	SQL	INTEL	2350
00000	33000000	0.85538104780917						
SQL	INTERNET	2015-11-22 20:17:39	70300000	915-8	2015-11-22 20:18:31	SQL	INTERNET	2350
00000	2147483647	0.26994665402126						
SQL	LOO	2015-11-22 20:17:39	34800000	915-8	2015-11-22 20:18:33	SQL	LOO	2350
00000	33600000	-8.26245979261802						
SQL	LTE	2015-11-22 21:00:10	9250000	915-2	2015-11-22 21:00:14	SQL	LTE	2350
00000	152000000	-8.11515651359090						
SQL	MOOC	2015-11-22 21:00:10	1680000	915-2	2015-11-22 21:00:15	SQL	MOOC	2350
00000	9650000	-1.5367813862359						
SQL	OPENSTACK	2015-11-22 21:00:11	9040000	915-2	2015-11-22 21:00:16	SQL	OPENSTACK	2350
00000	8950000	-8.80995216937006						
SQL	OTT	2015-11-22 21:00:11	4130000	915-2	2015-11-22 21:00:18	SQL	OTT	2350
00000	132000000	-8.19374918735121						
SQL	PHANTEC	2015-11-22 21:00:11	3820000	915-2	2015-11-22 21:00:19	SQL	PHANTEC	2350
00000	3040000	-1.632969180809						
SQL	PROCESS	2015-11-22 21:00:11	60300000	915-2	2015-11-22 21:00:20	SQL	PROCESS	2350
00000	177000000	0.2536428581235						
SQL	SAMSUNG	2015-11-22 21:00:11	3340000	915-2	2015-11-22 21:00:21	SQL	SAMSUNG	2350
00000	467000000	0.1352622386682						

Fig. 7 Result of Searching

그림 7은 검색이 끝난 후, 데이터베이스 테이블에 저장된 각 속성 값의 예를 보여주며, 그림 8은 사용자에게 보여주는 형태의 예이다.

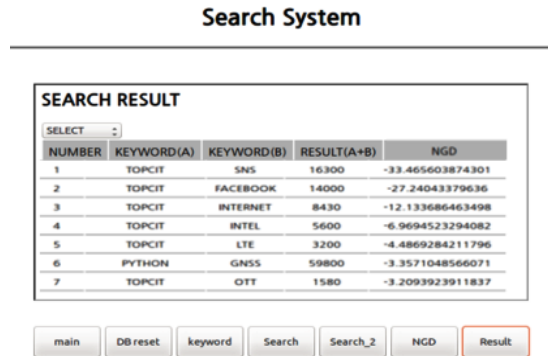


Fig. 8 End Result Screen

V. 결론 및 향후 연구

최근 많은 분야에서 연구하고 있는 빅데이터는 기존의 처리방법으로는 데이터를 분석하는데 많은 문제가 있어, 이에 맞는 새로운 분석 방법들이 연구되고 있다. 인터넷 역시 빅데이터의 대상으로 볼 수 있으며, 인터넷 상의 빅데이터를 분석하는 가장 쉬운 접근법 중 하나가 구글의 검색 결과 수를 활용하여 전체적인 추세를 분석하는 것이다.

구글의 검색 결과 수를 활용하여 추세를 파악한 예로는 미국의 독감환자 통계, 미 대선 예측, 한국 대선 예측, 뉴욕의 연방준비은행의 미국의 주택지표, 달러, 위안, 환율, 독일 실업률과 경제지표 예측을 들 수 있다.

만일 대선 후보 예측과 같이 특정 시점에서 특정 키워드를 검색한 결과의 수가 필요한 경우에는 비교적 간단한 검색만으로도 이러한 정보를 제공할 수 있겠지만, 보다 많은 키워드를 검색하거나 키워드 쌍을 조합하여 검색해야 하는 경우에는 이를 위한 전용 시스템을 사용하는 것이 보다 효율적이다.

본 논문에서는 복수개의 키워드를 동시에 검색하여 그 검색결과를 자동으로 추출하고 이를 관리할 수 있는 할 수 있는 구글 검색엔진을 활용한 키워드 검색결과 수 관리 시스템을 설계하고 구현하였다. 본 논문에서 제안하는 시스템은 크게 검색 에이전트, 저장 노드, 검

색 노드로 구성되어 동작한다. 최종 검색 결과로는 검색 키워드, 검색 결과 수, 검색 결과 수를 활용하여 두 키워드의 거리를 계산하는 NGD(Normalized Google Distance)가 제공된다.

본 논문에서 제안하는 시스템을 활용하여 임의의 키워드 집합에서 서로 관련이 깊은 두 키워드 쌍을 쉽게 찾을 수 있고, 또한 이러한 대용량 키워드 쌍의 집합들에 대한 시간에 따른 관련성의 변화 등을 파악할 수 있는 시작점이 될 수 있다.

향후 연구 과제로는 제안 시스템의 확장성에 관한 부분과 실제 데이터로 적용시켜 활용하는 부분 등이 있다.

REFERENCES

- [1] S. Brin and L. Page, "The Anatomy of a large-scale hypertextual Web search engine", *Computer Networks and ISDN Systems*, vol.30, no.1-7, pp.107-117, Apr. 1998.
- [2] Google Flu Trend [Internet]. Available: <https://www.google.org/flutrends/about/>
- [3] R. L. Cilibrasi and P. M. Vitanyi, "The google similarity distance", *IEEE Transactions on, Knowledge and Data Engineering*, vol. 19, no. 3, pp.370-383, Mar. 2007.
- [4] Google Trend [Internet]. Available: <https://www.google.com/trends>
- [5] J. Ginsberg, M. H. Mohebbi, R. S. Patel, L. Brammer, M. S. Smolinski and L. Brilliant, "Detecting influenza epidemics using search engine query data," *Nature*, vol. 457, pp. 1012-1014, Feb. 2009.
- [6] H. Achrekar, A. Gandhe, R. Lazarus, S. H. Yu, and B. Liu, "Predicting flu trends using Twitter data," *The First International Workshop on Cyber-Physical Networking Systems*, pp.702-707, Apr. 2011.
- [7] Google Ngram Viewer [Internet]. Available: <https://books.google.com/ngrams>
- [8] F. H. Messerli, "Chocolate Consumption, Cognitive Function, and Nobel Laureates," *The New England And Journal Of Medicine*, vol. 367, pp. 1562-1564, Oct. 2012.



이주연(Ju-Yeon Lee)

2013년 동의대학교 문헌정보학과 문헌정보학사
 2015년 동의대학교 컴퓨터소프트웨어공학과 공학석사
 ※관심분야 : 데이터베이스, 정보검색



이중화(Junghwa Lee)

1992년 부산대학교 전자계산학과 학사
 1995년 부산대학교 전자계산학과 석사
 2001년 부산대학교 전자계산학과 박사
 2002년 ~ 현재 동의대학교 컴퓨터소프트웨어공학과 교수
 ※관심분야 : 데이터베이스, 한글정보처리, 멀티미디어시스템



박유현(Yoo-Hyun Park)

1996, 1998, 2008년 부산대학교 전자계산학과 이학사, 이학석사, 이학박사
 2000년 한국국방연구원(KIDA) 연구원
 2001년 ~ 2009년 한국전자통신연구원(ETRI) 선임연구원
 2009년 ~ 현재 동의대학교 컴퓨터소프트웨어공학과 부교수
 2012년 ~ 2014년 동의대학교 부산T융합부품연구소 부소장
 ※관심분야 : 인터넷 시스템, 클라우드 시스템, 빅데이터, 소프트웨어 품질, IT 융합 서비스