

---

# 쉽게 접하는 빅데이터 알고리즘

---

오수진 주임

# 목차

## 1. BigData Life

## 2. 우리가 관심있어 하는것

### 2-1. 실생활 예

## 3. 분석과정

### 3-1. 뉴스의 본문 추출

### 3-2. 형태소 분석

### 3-2. TF-IDF 알고리즘

## 4. 도출된 결론

### 4-1. 분야별 주요 키워드

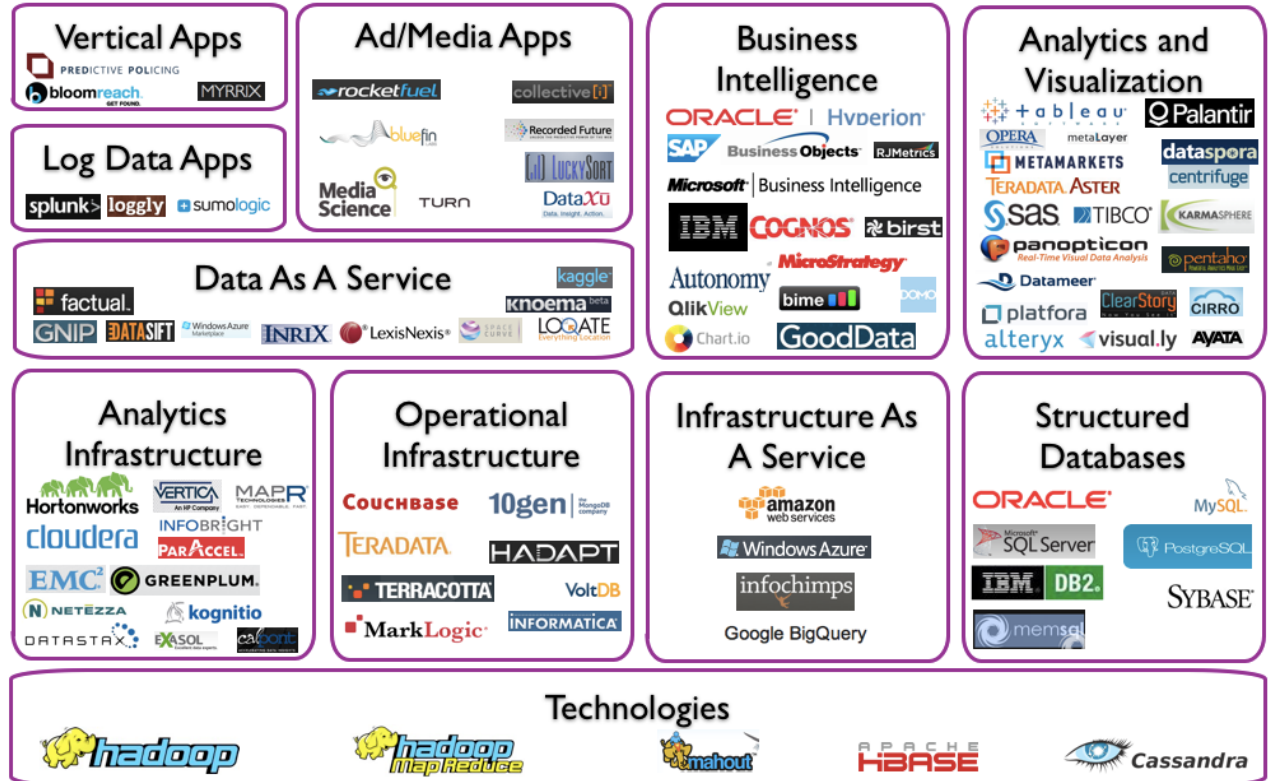
## 1. BigData Life



빅데이터라는 기술이 우리의 삶과 얼마나 밀접한 관련이 있을까?

## 1. BigData Life

# Big Data Landscape



Copyright © 2012 Dave Feinleib

[dave@vcdave.com](mailto:dave@vcdave.com)

[blogs.forbes.com/davefeinleib](http://blogs.forbes.com/davefeinleib)

빅데이터를 분석하는데에는 수많은 기술들이 있습니다.



# Big Data Landscape (Version 2.0)

## Infrastructure



## Analytics



## Applications



## Data Sources



## Open Source Projects





# Big Data Landscape (Version 2.0)

## Infrastructure



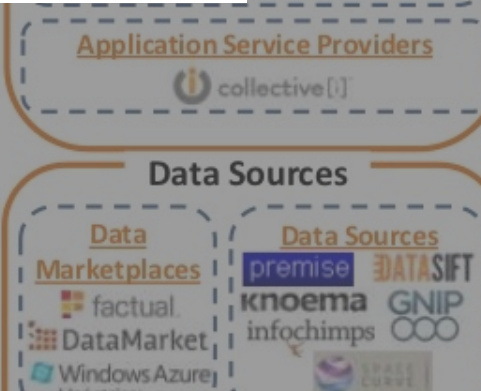
## Analytics



## Applications



하... 빅데이터를 분석하는 많은 기술들이  
도대체 어떻게 우리의 삶과 관련이 있다는거지...??



## Open Source Projects





# Big Data Landscape (Version 2.0)

## Infrastructure

**NoSQL Databases**  
10gen, DATASTAX, basho, COUCHBASE, CLOUDANT, HYPERTABLE, Neo4j, SCORIS, Oracle Database  
**NewSQL Databases**  
MarkLogic, paradigm4, memsql, SQLFire, DRAWNSCALE, VoltDB, NUODB

**MPP Databases**  
VERTICA, An HP Company, Kognitio, PARACCEL, GREENPLUM, A DIVISION OF EMC, TERADATA, N, NETEZZA, InfiniDB, Microsoft SQL Server  
**Storage**  
Cleversafe, panasas, nimblestorage, AMPLIDATA, Compuverde

**Management / Monitoring**  
OUTER, THOUGHT, oceanSYNc, StackIQ, bundy, DATADOG  
**Crowdsourcing**  
CROWD COMPUTING systems, CrowdFlower

**Hadoop Related**  
cloudera, HADAPT, Hortonworks, infochimps, MAPR, HSTREAMING, Zettaset, MORTAR, Microsoft, GREENPLUM, A DIVISION OF EMC, amazon, Qu, bole, agril

**Cluster Services**  
LexisNexis, HPCC Systems, Acunote, Secur, Storm, iMPE, TRACE/VE, codefortytw, DATAGUISE  
**Collection /**

## Analytics

**Analytics Solutions**  
Palantir, platforma, PERVASIVE, Datameer, KARMASHERE, DataHive, DIGITAL REASONING, dataspora, PRECOG

**Statistical Computing**

**Data Visualization**  
Quid, visual.ly, ACTUATE, Kitenga, metaLayer, ClearStory, +tableau, ISS, Quantum4D

**Social Media**  
simple reach, Dataminr, ops Services, BIG, Company, OPERA, search, Autonomy, T Analytics, sumologic, sourced, SMB Analytics, sumall

## Applications

**Ad Optimization**  
DataXu, aggregate knowledge, m6d, MediaMath, bluekai, ai Match, rocketfuel, thetradedesk, TURN, 33 across  
**Publisher Tools**  
VISUAL, revenue, Yielddex, yieldbot  
**Marketing**  
LATTICE ENGINES, Sailthru, SCIENCE, bloomreach, CLICKFOX

**Industry Applications**  
NEXT BIG SOUND, KNEWTON, rest, cash, wonga, numberFire, Mile Sense, BILLI, Climate Solutions, Bloomberg, GUARD  
**Application Service Providers**  
collective[i]

## Data Sources

**Data Marketplaces**  
premise, DATA SIFT, knoema, Gnip, infochimps, OOO, BASIS, fitbit

## Framework



하둡이란 기술을 실생활에서 어떻게 사용 할 수 있는지 알아보도록 하겠습니다.

## Open Source Projects

**Framework**  
Hadoop MapReduce, HDFS  
**Query / Data Flow**  
hive, flume  
**Database Access**  
Cassandra, SciDB, HBASE, CouchDB, Sqoop  
**Coordination / Workflow**  
ZooKeeper, talend  
**Real-Time**  
Storm  
**Statistical Tools**  
SciPy  
**Machine Learning**  
mahout  
**Cloud Deployment**

## 2. 우리가 관심있어 하는것

분명 나는 이것이 필요하다는 것을  
알지만 과연 나의 친구, 가족들도  
이런 기술들이 필요하다는 것을

“ 공감 ”

할 수 있을까요?

[ NEWS ]

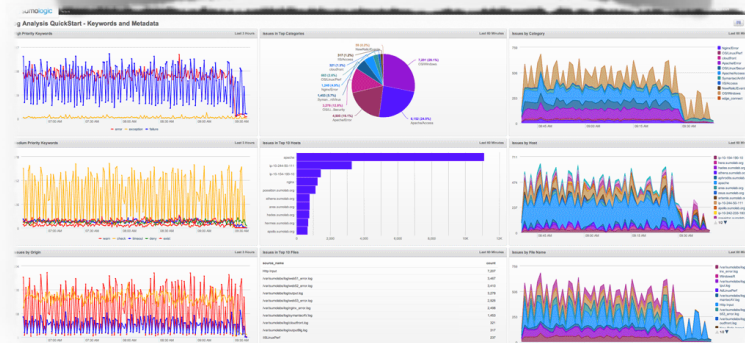
### SK C&C, 빅데이터를 활용한 실시간(스트리밍) 통합보안로그분석 플랫폼 개발

1,800여개 사이트 대상 로그 데이터 분석으로 해킹 탐지·검색·추적까지 한번에 처리  
실시간 로그 데이터 분석으로, 해킹 탐지시간 실시간(스트리밍) 데이터 처리 가능  
IDS-방화벽 통합 분석 통해 해킹 이상 징후까지 해킹 탐지 정확도 향상  
3TB의 빅데이터 대상 해킹 검색, 4초 이내로 단축해 국내 최고 수준 평가  
공개 S/W 활용으로 향후 사용 고객의 비용 절감 효과 기대

### Daum 실시간 분석 사례

- 로그 분석 사례**
  - 전사 로그를 통한 통계 분석
  - 광고 로그 분석을 통한 타겟팅
  - 검색 품질 향상 분석 및 개선
  - 광고 및 클릭 로그 분석을 통한 타겟팅
  - 카페 로그 분석을 통한 사용자 카페 추천
  - 게임 서버 로그 분석 등
- 데이터 분석 사례**
  - 쇼핑 마우스 상품 클릭 분석 사례
  - 다음 Top 토픽 분석 및 추천 서비스
  - UCC 문서의 소셜 유저 필터링
  - 사물 검색 이미지 역색인
  - 자연어 처리 텍스트 분석
  - 모바일 광고 데이터별 매체 분석 등
- 연구 개발 사례**
  - 이미지 유사성 매칭 분석
  - 대용량 시맨틱 웹 검색 엔진 개발
- 서비스 적용 (MongoDB/카산드라)**
  - 마이 아고라
  - 검색 광고 노출 최적화
  - 최근 방문 카페 저장
  - 사내 캐시 서버(Redis)
  - 사내 Git 저장소(Redis)
- 데이터 처리 (Hbase)**
  - 검색 엔진 색인 문서 저장
  - 서버 모니터링 데이터 저장
  - 로그인 로그 저장
  - 카페 방문 로그 저장
- 서비스 분석**
  - 미디어 다음 실시간 분석
  - 모바일/PC탑 실시간 분석
- 데이터 수집**
  - Twitter 실시간 데이터 수집기

real time





## 2. 우리가 관심있어 하는것



대한민국 국민이 매일매일 빼먹지 않고 하는 필수 코스

➡ **NAVER** 뉴스보기

## 2. 우리가 관심있어 하는것

그런데 뉴스를 보다가 한가지 궁금증이 생겼습니다.

최신기사 정치 경제 국제 사회 문화 의학과학 사람속으로

이렇게 다양한 분야에서



이렇게 다양한 뉴스가 있는데

과연 분야별로 어떤 키워드가 가장 관심을 받을까요??

### 3. 분석하기

클라라 거짓말 모음 '끝도 없네'



#### 본문

클라라의 거짓말이 새삼 화제다.

28일 방송된 SBS '한밤의 TV연예'에서는 스타들을 향한 대중의 분노와 스타들의 잘못된 대처에 대한 분석을 전했다.

최근 논란의 중심에 선 클라라는 현재 폴라리스엔터테인먼트와 갈등을 빚고 있다. 클라라는 소속사 회장으로부터 성적 수치심을 느꼈다고 주장했지만 한 매체의 문자 공개 후 여론은 급변했다.

각 기사의 본문만 발췌하여 키워드 별로 나누고 이를 분석하여 어떤 키워드가 **HOT** 한지 알아보고 싶었습니다.



## 3. 분석하기

1. 가장 먼저 rss와 xml 파서를 이용해 분야별로 기사의 링크를 수집했습니다.

최신기사 정치 경제 국제 사회 문화 의학과학 사람속으로

```

▼<item>
  <title>클라라, 그간 했던 거짓말들 모아보니...누리꾼 '분노'</title>
  ▼<link>
    http://mbn.mk.co.kr/pages/news/newsView.php?category=mbn00012&news\_seq\_no=2189396
  </link>

```

2. 기사는 대부분 “~다.” 라고 끝내는 평서문의 형태로 이루어져 있습니다.



클라라의 거짓말이 새삼 화제다.

28일 방송된 SBS '한밤의 TV연예'에서는 스타들을 향한 대중의 분노와 스타들의 잘못된 대처에 대한 분석을 전했다.

최근 논란의 중심에 선 클라라는 현재 폴라리스엔터테인먼트와 갈등을 빚고 있다. 클라라는 소속사 회장으로 부터 성적 수치심을 느꼈다고 주장했지만 한 매체의 문자 공개 후 여론은 급변했다.

## 3. 분석하기

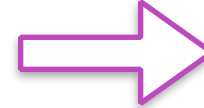
### 3. 본문만 가져오기



클라라의 거짓말이 새삼 화제다.

28일 방송된 SBS '한밤의 TV연예'에서는 스타들을 향한 대중의 분노와 스타들의 잘못된 대처에 대한 분석을 전했다.

최근 논란의 중심에 선 클라라는 현재 폴라리스엔터테인먼트와 갈등을 빚고 있다. 클라라는 소속사 회장으로 부터 성적 수치심을 느꼈다고 주장했지만 한 매체의 문자 공개 후 여론은 급변했다.

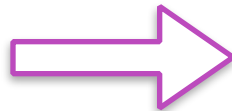


클라라의 거짓말이 새삼 화제다.  
28일 방송된 SBS '한밤의 TV연예'에서는 스타들을  
향한 대중의 분노와 스타들의 잘  
못된 대처에 대한 분석을 전했다.

### 4. 키워드 추출하기 - 형태소 분석

클라라Noun, 의Josa, 거짓말Noun, 이Josa, 새삼Noun, 화제Noun,  
다Josa, .Punctuation, 28Number

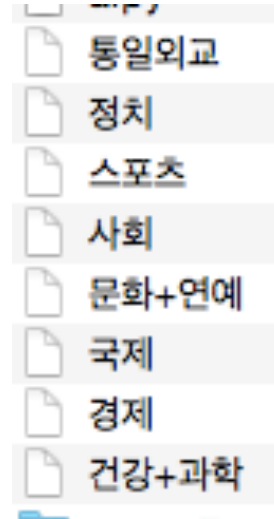
기사 본문의 명사만 잘라내어 분류 합니다.



클라라  
나이  
실제  
나이  
살  
한국

### 3. 분석하기

5. 분야별 기사의 키워드를 TXT 파일로 저장합니다.



6. 수집된 데이터를 TF-IDF 알고리즘을 이용하여 분석합니다.

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

**TF-IDF**

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents



## 3. 분석하기

### 1. Term Frequency : tf

단어 빈도수는 말 그대로 "단어가 그 문서에서 나타난 횟수"를 나타냅니다.

	문화 연예	사회	스포츠	정치	국제	통일 외교	경제
클라라	110	4	7	0	0	0	0
제일제당	0	0	0	0	0	0	16
한국	31	32	3	1	21	1	11

문서 내에서 각 단어의 중요도 파악 !!!

### 3. 분석하기

#### 1. Term Frequency : tf

단어 빈도수는 말 그대로 "단어가 그 문서에서 나타난 횟수"를 나타냅니다.

	문화 연예	사회	스포츠	정치	국제	통일 외교	경제
클라라	110	4	7	0	0	0	0
제일제당	0	0	0	0	0	0	16
한국	31	32	3	1	21	1	11

문서 내에서 각 단어의 중요도 파악 !!!

### 3. 분석하기

#### 2. Document Frequency : df

문서 빈도수는 "해당 단어가 나타난 문서의 수" 입니다.

내가 찾을 문서에선 많이 나타나고, 다른 문서에선 적게 나타날 수록 중요하다.

	문화 연예	사회	스포츠	정치	국제	통일 외교	경제
클라라	110	4	7	0	0	0	0
제일제당	0	0	0	0	0	0	16
한국	31	32	3	1	21	1	11

위 차트에서 제일제당의 df 값은 1이고 한국의 df 값은 7입니다.



## 3. 분석하기

### 2. Document Frequency : df

문서 빈도수는 "해당 단어가 나타난 문서의 수" 입니다.

내가 찾을 문서에선 많이 나타나고, 다른 문서에선 적게 나타날 수록 중요하다.

	문화 연예	사회	스포츠	정치	국제	통일 외교	경제
클라라	110	4	7	0	0	0	0
제일제당	0	0	0	0	0	0	16
한국	31	32	3	1	21	1	11

위 차트에서 제일제당의 df 값은 1이고 한국의 df 값은 7입니다.

### 3. 분석하기

### 3. Inverse Document Frequency : idf

역문서 빈도수는 df를 역수 취한 것입니다.

df값이 클 수록 중요하지 않은 단어를 나타내는 것인데,

이것을 반대로 값이 클수록 중요한 단어로 나타내기 위하여 역수를 취하는것입니다.

	tf	df	전체 문서의 개수(N)	idf
제일제당	<b>16</b>	<b>1</b>	<b>7</b>	<b>7/1</b>

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

**TF-IDF**

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents

## 4. 도출된 결론

	tf	df	전체 문서의 개수(N)	idf
제일제당	<b>16</b>	<b>1</b>	<b>7</b>	<b>7/1</b>

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

### TF-IDF

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents

제일제당의 TF-IDF값 :  $16 * \log(7/1) = 35.155$

## 4. 도출된 결론

제일제당의 TF-IDF값 :  $16 * \log(7/1) = 35.155$

경제	오닝	35.156
경제	제일제당	35.156
경제	소상	22.181
경제	아미노산	21.972
경제	제품	21.972
경제	창업	21.972
경제	한전	20.794

각 키워드의 tf-idf값을 기준으로 정렬을 하면  
분야별 핵심 키워드를 찾을 수 있습니다.



## 4. 도출된 결론

제일제당의 TF-IDF값 :  $16 * \log(7/1) = 35.155$ 

경제	오닝	35.156
경제	제일제당	35.156
경제	소싱	22.181
경제	아미노산	21.972



**CJ제일제당** 메치오닌시장 본격 공략 파이낸셜뉴스 | 22분전 | 네이버뉴스 | [🔗](#)

말련 메치오닌공장 가동 "5兆 글로벌시장 1위 도약" CJ제일제당이 50억달러(5조4310억원) 규모의 사료용 아미노산 메치오닌 시장에 본격 진출한다. 29일 CJ제일제당은 프랑스 아르케마사와 손잡고 말레이시아에 총 4억달러...

↳ CJ제일제당, 5조원 규모 메치오닌 시장 본격 공략 아시아경제 | 9시간전 | 네이버뉴스

↳ CJ제일제당, 메치오닌 시장 공략 강화..세계 1위... 이데일리 | 9시간전 | 네이버뉴스

[관련뉴스 전체보기 >](#)

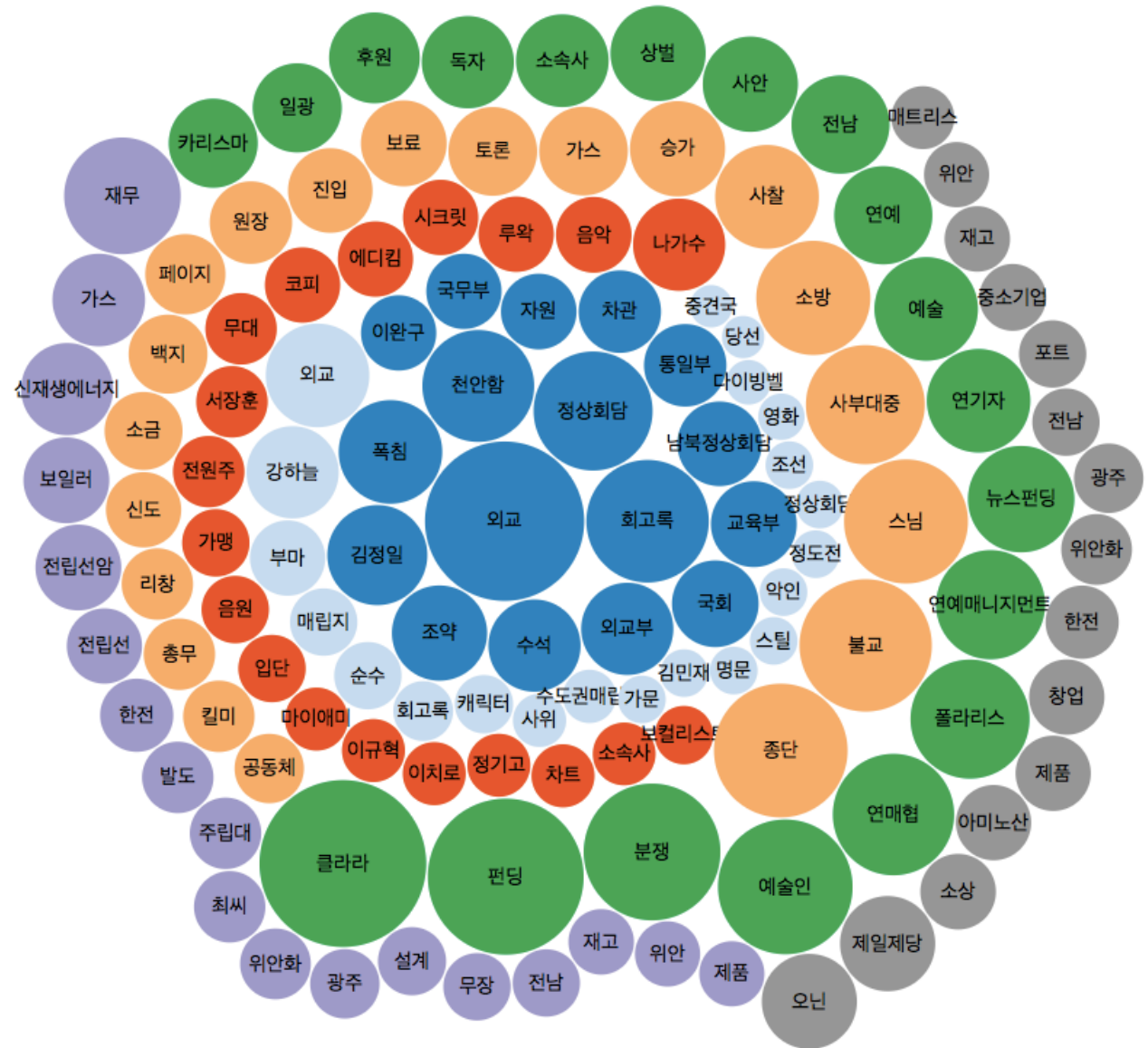
[관련뉴스 48건 전체보기 >](#)

[관련뉴스 18건 전체보기 >](#)

메치오닌의 형태소 분석결과 : [메Verb, 치PreEomi, 오닝Noun]

**(D3 bubble chart)**

**(D3 bubble chart)**





미래를 예측하는 가장 좋은 방법은 미래를 창조하는 것이다.

— 피터 드러커 —

---

# Q & A