

Big Data Near You

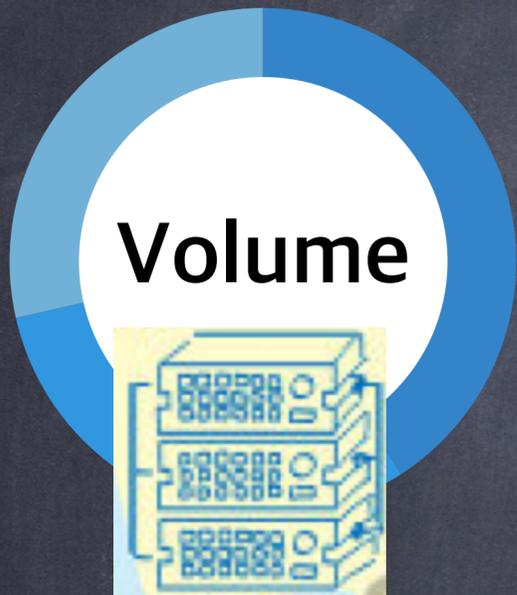
2016. 06. 16

Prof. Sehyug Kwon



Dept. of Statistics

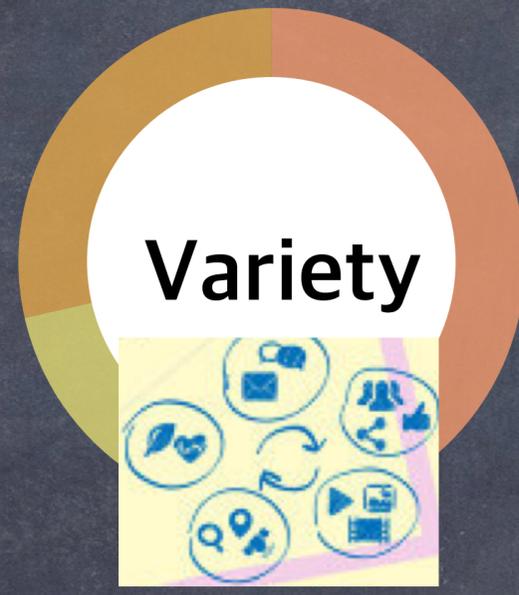
4V's of Big Data



Volume

대용량

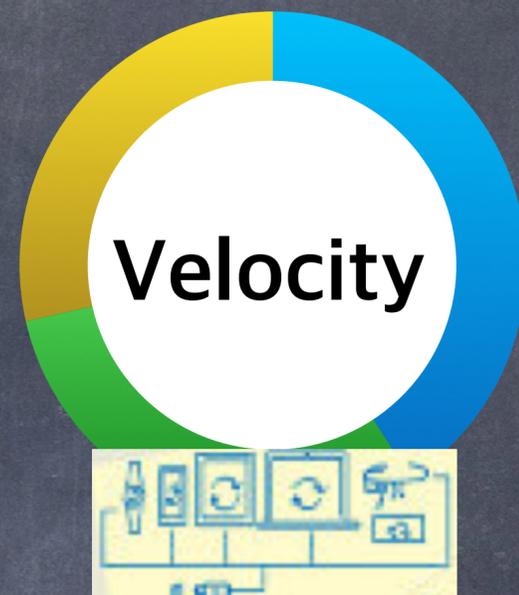
- 1 exa(1,000⁴) - peta -exa-zetta(1000⁷) bytes in 2020
- 통화내역, 카드사용내역
- 회사, 정부: DB-DW-Data mart



Variety

다양한 유형

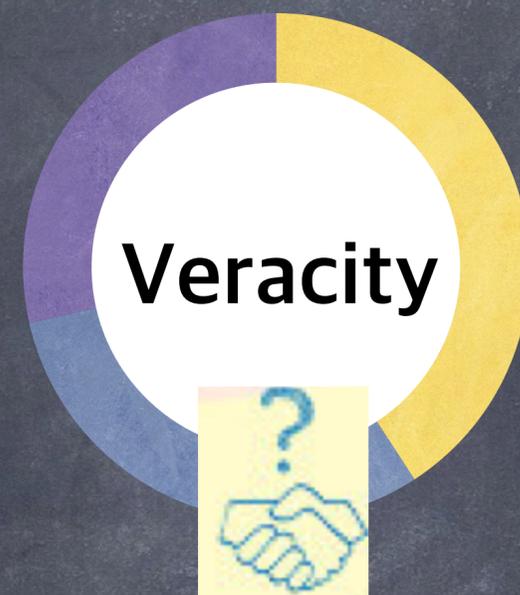
- 텍스트 마이닝
- 이미지, 멀티미디어
- (비)정형화 포맷



Velocity

실시간 정보

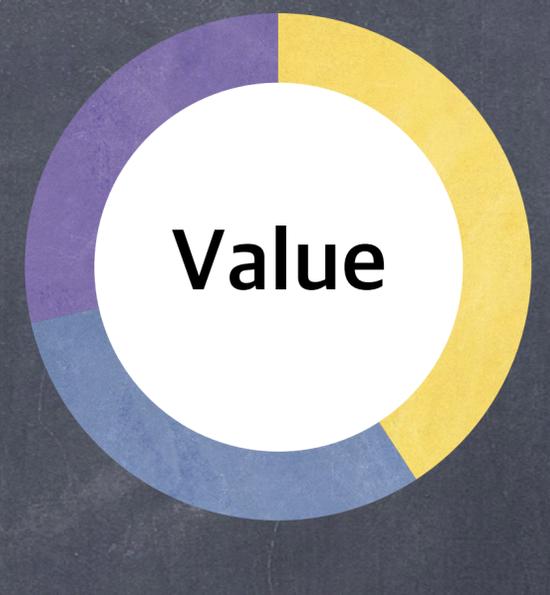
- Streaming data
- 실시간 분석결과



Veracity

(불)확실성

- 데이터 이력
- 불확실성=비용
- 1/3 CEO 의사결정 불신



Value

가치

- 딥러닝, 학습효과
- 비즈니스 정보

Three Experts in Big Data

- Hadoop (분산파일처리)
- MapReduce (분산프로그래밍모델)
- Java / Python / Ruby
- NoSQL, DB
- Apache Spark

개발자

관리자

- 하둡, 리눅스 관리
- Cluster Management
- Cluster Performance
- Virtualization

- 데이터 과학 - 모델링
- 기계학습, 마이닝기법
- 빅데이터 벤더 : R/SAS
- 데이터 Visualization

데이터
분석가

Leaders in Big Data

Deep Learning Tools

New this year was a category of Deep Learning

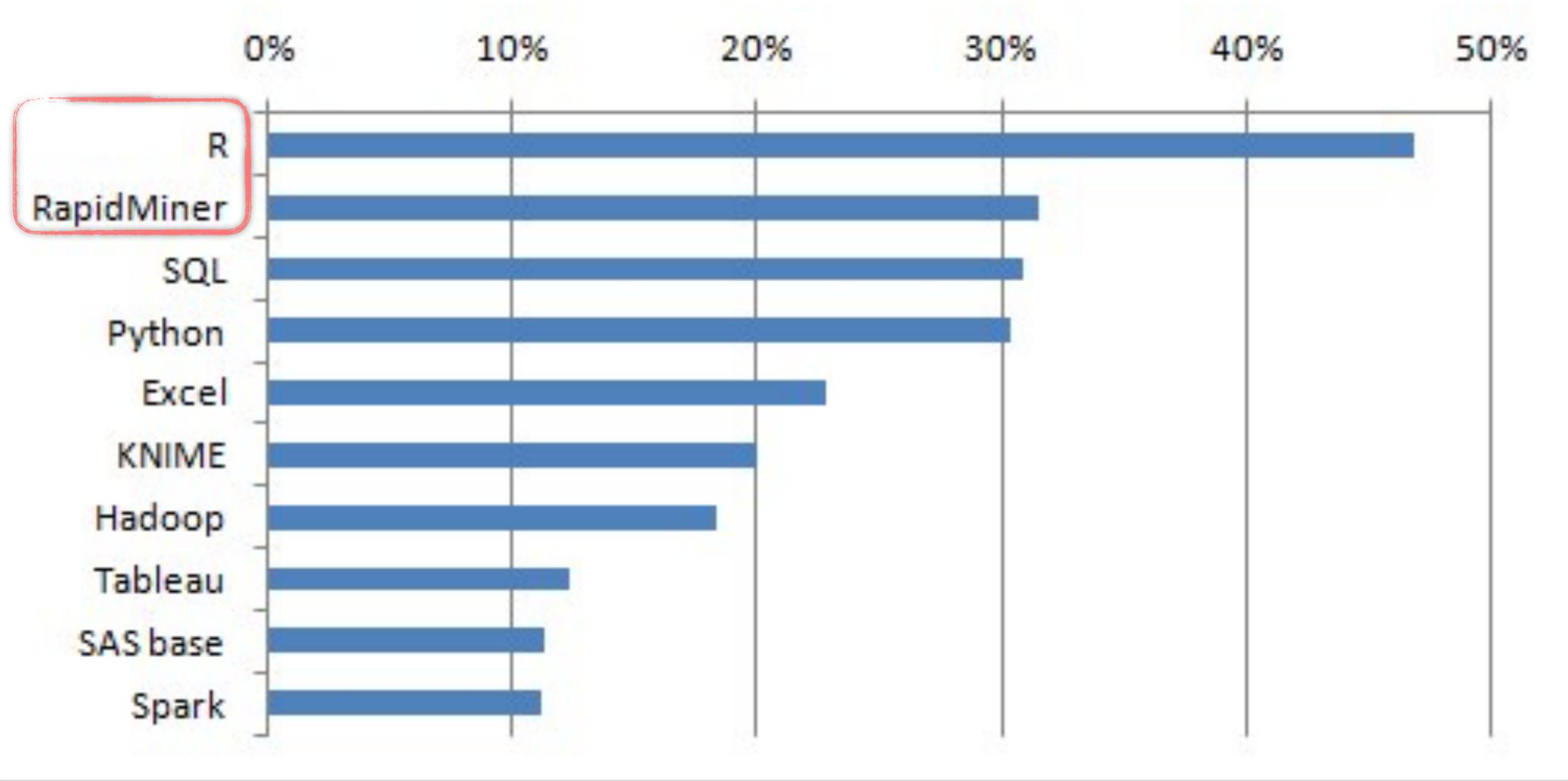
- Pylearn2 (55 users)
- Theano (50)
- Caffe (29)
- Cuda-convnet (17)
- Deeplearning4j (12)
- Torch (27)

Top Hadoop/Big Data tools were

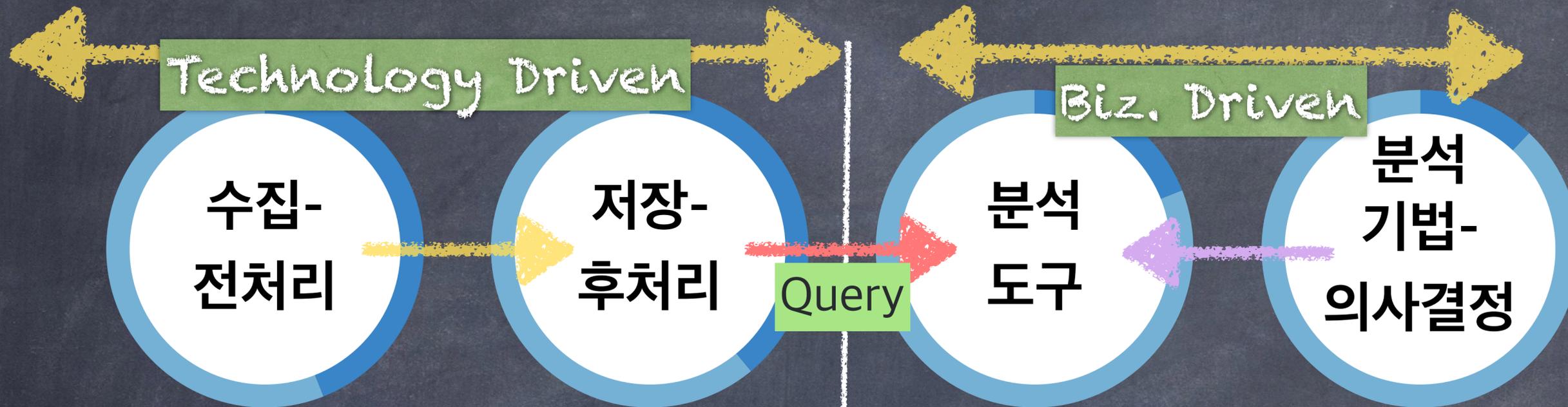
- Hadoop, 18.4% share (507 votes)
- Spark, 11.3% (311)
- Hive, 10.2% (282)
- SQL on Hadoop tools, 7.2% (198)
- Pig, 5.4% (150)
- HBase, 4.6% (127)
- Other Hadoop/HDFS-based tools, 4.5% (125)
- MLlib, 3.3% (91)
- Mahout, 2.8% (76)
- Datameer, 0.8% (23)

Leaders in Big Data (cont.)

Top Analytics, Data Mining, Data Science software used, 2015



Big Data Flow



Collecting

Processing

Analysis

Visualization



HDFS

Hive

Mahout

SAS-Insight

Warehousing + E-Minor

open source R

- Unsupervised Learning
- Social Media analytics
- Sentiment analysis
- Predictive modeling
- Visualization
- Simulation

Case in Big Data

map 거미줄 프로젝트

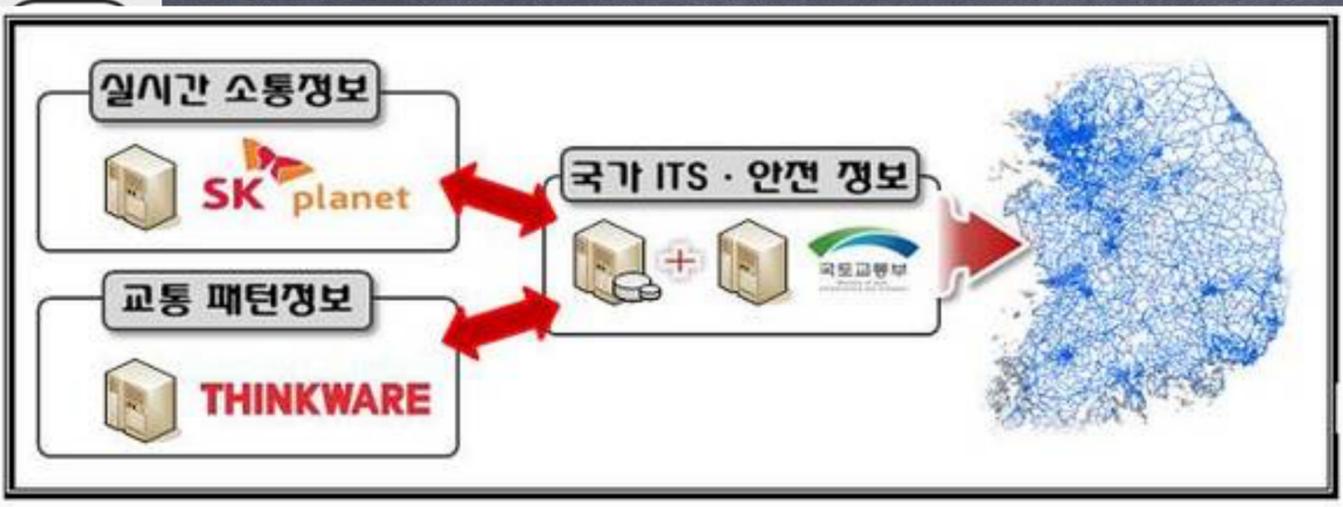
이성당

이성당에서 인기있는 다음 목적지는?

2위 히로쓰가옥

전라북도 군산시 신흥동

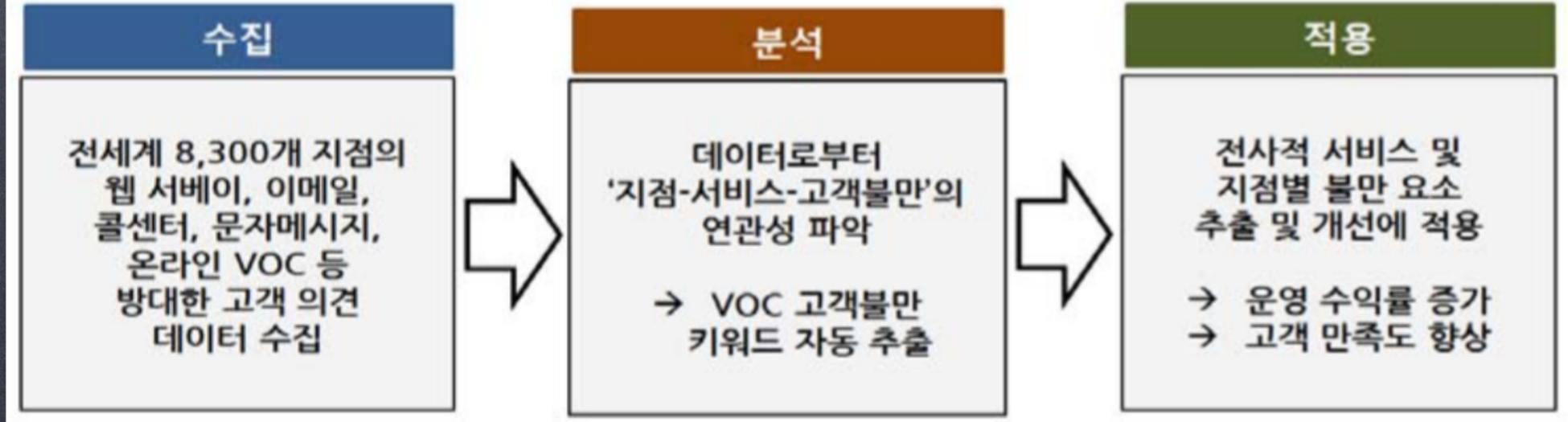
지도보기 map 검색



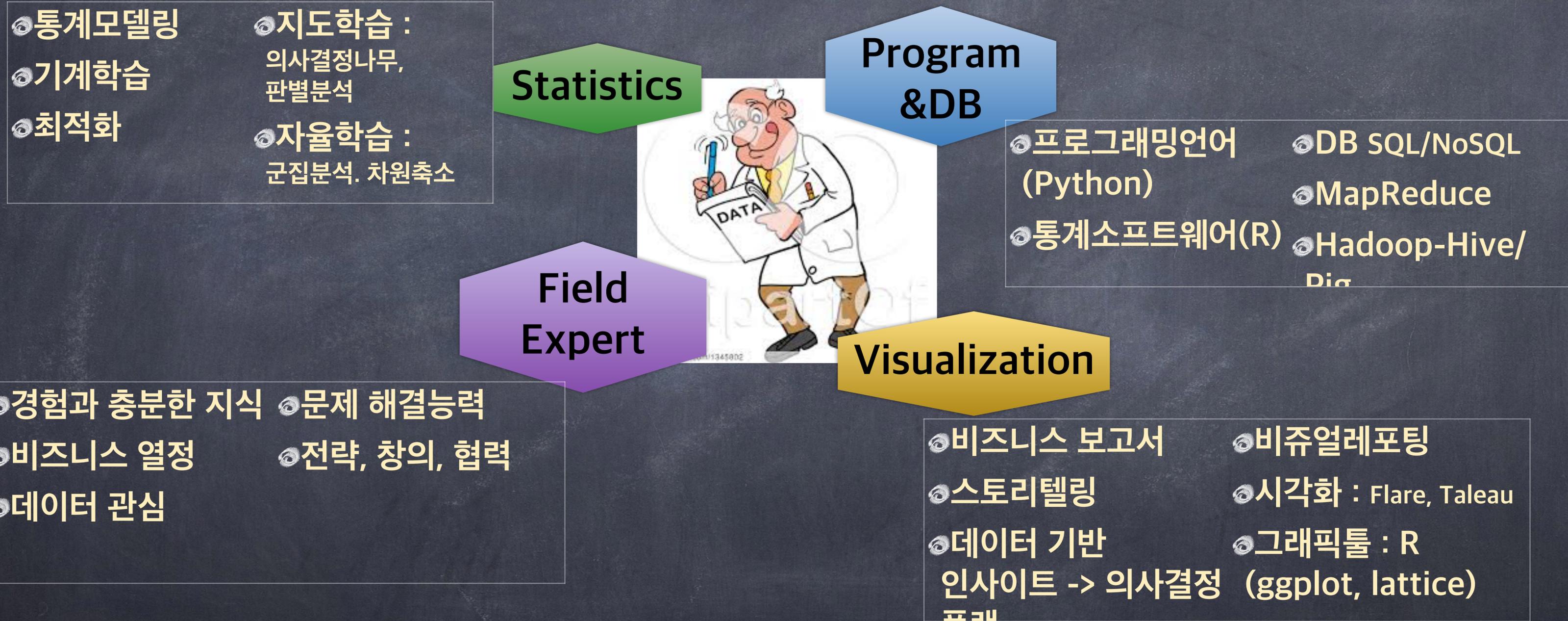
Kakao + =

KakaoTaxi

[그림] 허츠 빅데이터 적용 프로세스



Does this Data Scientist exist in the world?



Supervised Learning 기계학습

- 기계학습 기법에서는 최적 모형의 탐색 및 선정은 통계분석가에 의해 진행 - 예측 및 판별이 주 목적
- 훈련 Trained(목표변수(target)=함수(예측변수 predictors) 관계를 도출)-검증 Validation(모형의 타당성을 검증)-평가 Test(서로 다른 통계적 방법들의 평가)
- 인공지능 (컴퓨터 학습할 수 있도록 하는 알고리즘과 기술) 반복 학습 알고리즘 이용한 데이터 인사이트 탐색 - (예) 스팸 메일 검증여부
- Deep Learning: (**Unsupervised Learning**)
 - ▶ 기계 학습의 한 영역으로, 특히 음성/텍스트/이미지 인식 분야에서 획기적인 발전을 거듭하며 급성장하고 있음,
 - ▶ 기본 원리 : 다수의 히든 레이어를 갖춘 신경망을 통해 컴퓨터가 태스크를 학습하고 정보를 체계화하여 스스로 패턴을 찾아낼 수 있게 하는 것입니다.

Supervised Learning (cont.)

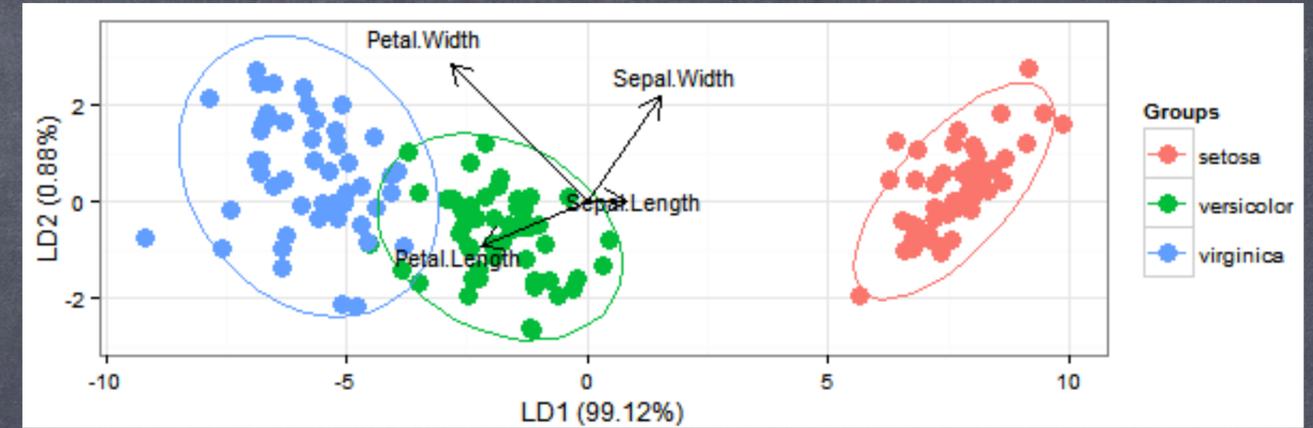
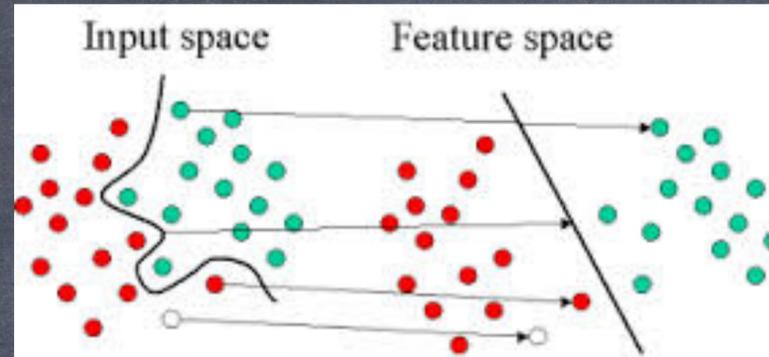
판별분석, 예측모형 : 신용평가 모형

▶ 예측변수(input) 활용 판별규칙

SVM

▶ 패턴인식, 지도학습

▶ 분류, 회귀분석에 사용



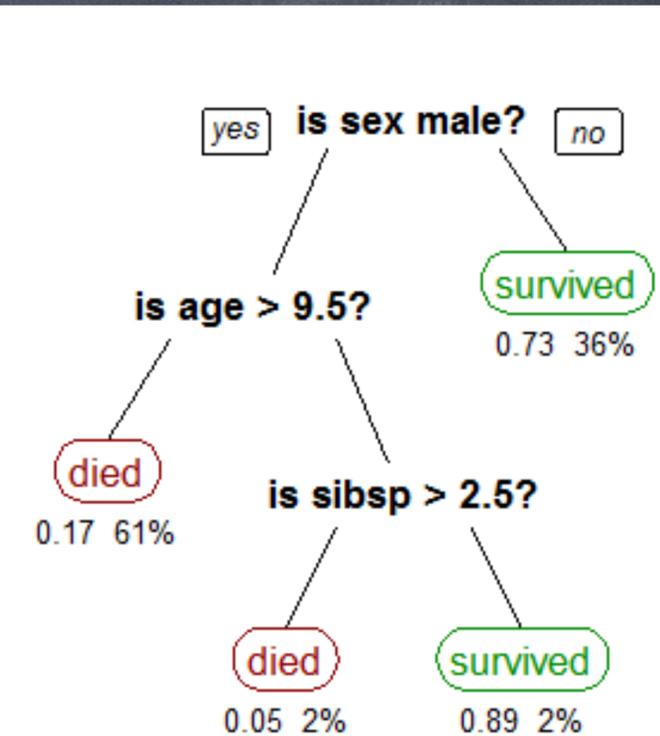
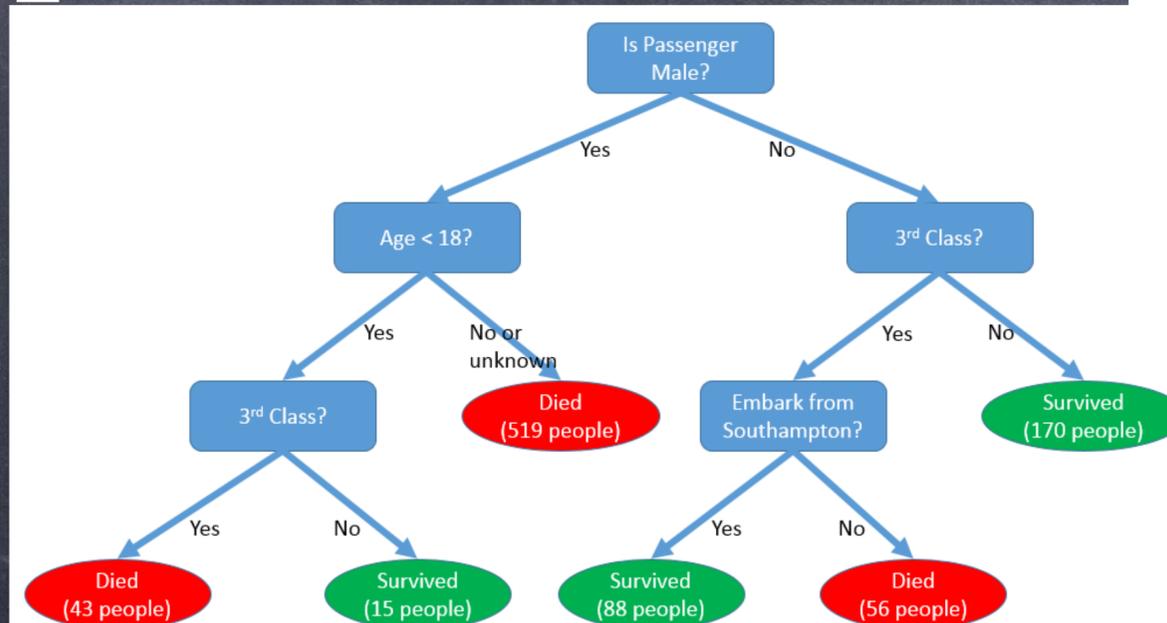
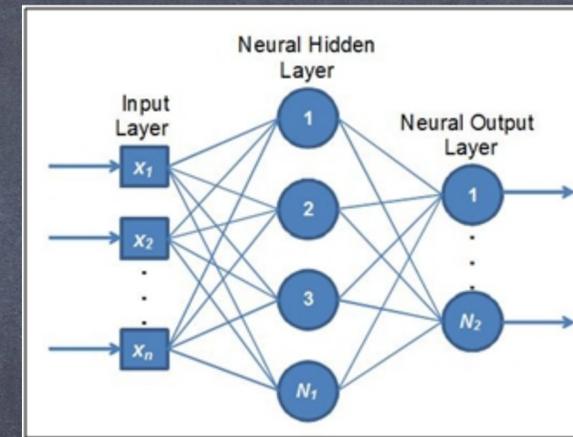
Artificial Neural Network 인공신경망

▶ 인공뉴런(노드)이 학습을 통한 문제해결

▶ 교사학습 vs. 비교사 학습(기계학습)

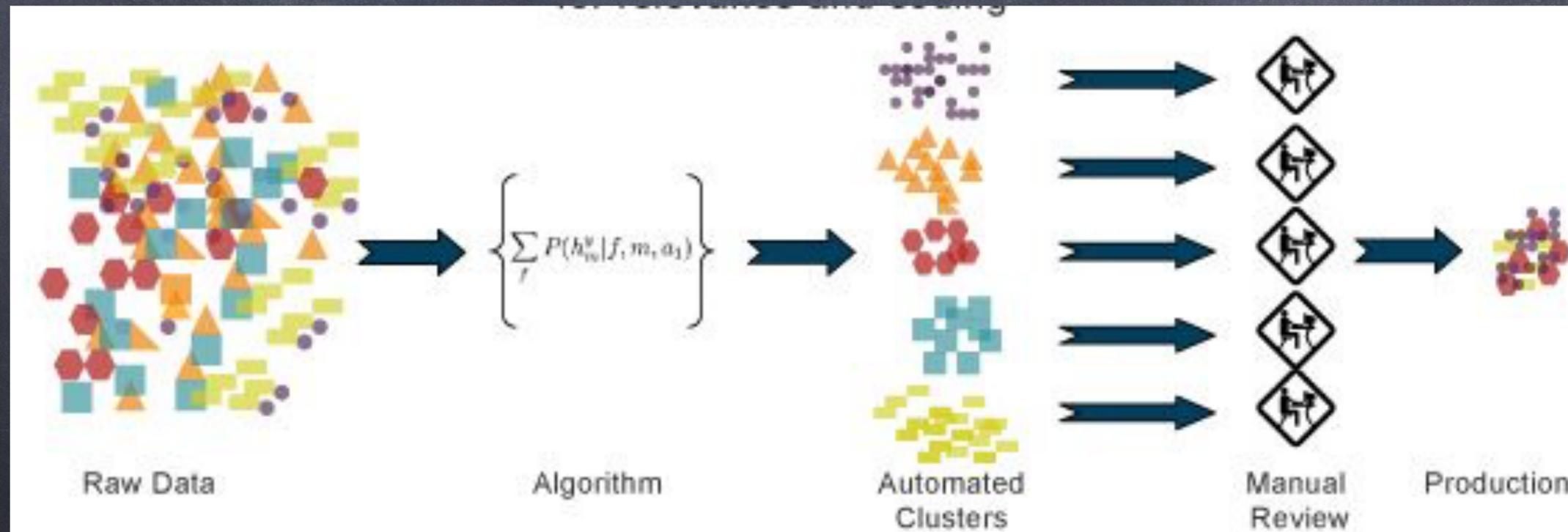
의사결정나무

▶ 타이타닉 생존자 분류



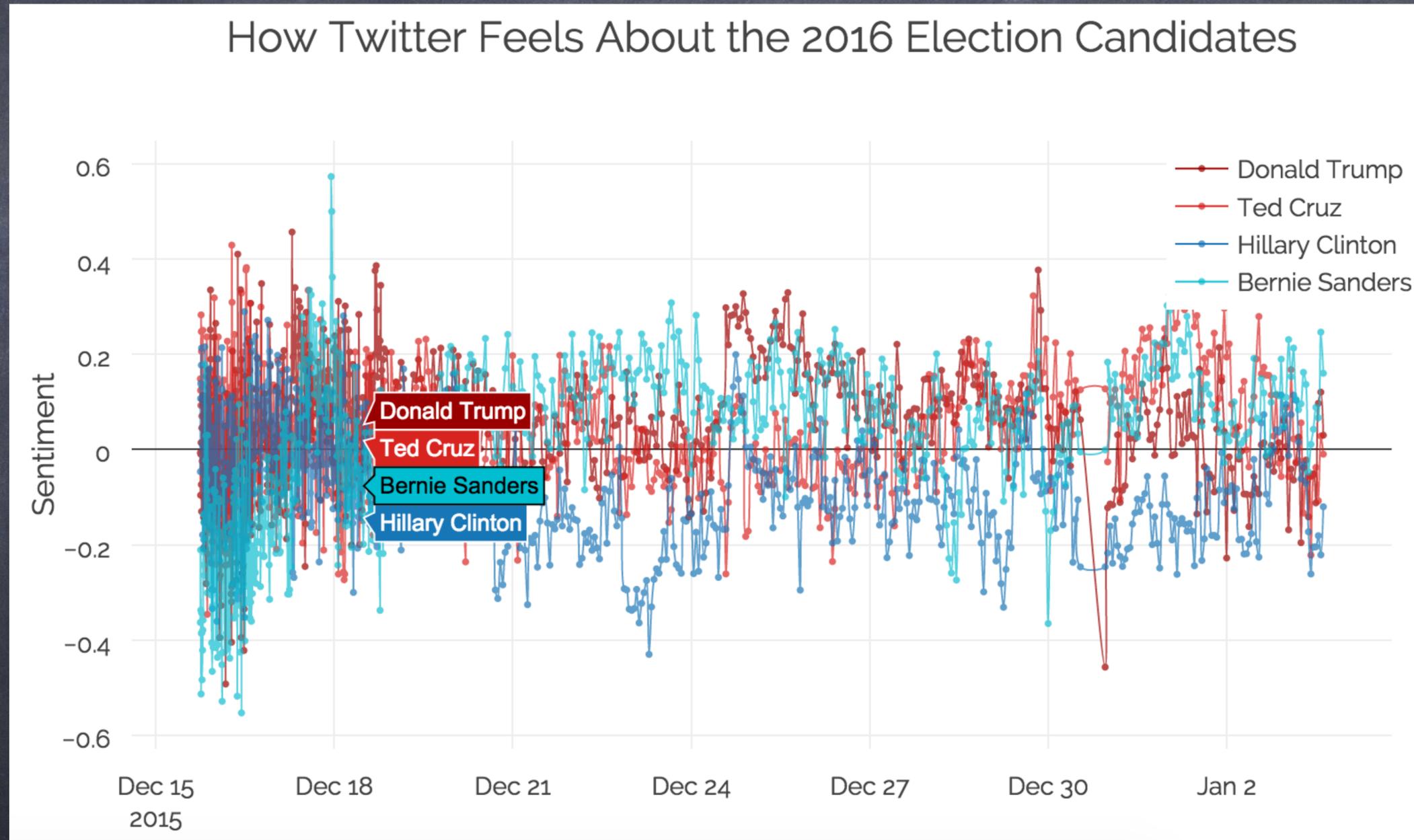
Un-Supervised Learning 자율학습

- 군집분석 : 집단 분류가 없는 개체들을 개체 내의 내재된 관계를 설명하는 함수(유사성의 함수)를 활용하여 분류
 - ▶ 계층적 군집 - (Linkage 연결-덴드로그램)
 - ▶ 비계층적 군집 - K-means
- 신경망 이론
 - ▶ SOM(self organizing map), ART(adaptive resonance theory) 알고리즘



Sentiment Analysis

Opinion Mining 오피니언 마이닝



Twitter mining with R

R-설치

<http://r-project.org>

```
1 #R=Twitter mining with R
2 Sys.getlocale() #set default
3 library(twitteR) #install.packages("twitteR")
4 library(RCurl) #install.packages("RCurl")
5 library(ROAuth) #install.packages("ROAuth")
6 library(base64enc) #install.packages("base64enc")
```

[[95]]

[1] "3RMarket0: 빅데이터의 분류\nhttps://t.co/Apjm3ybiwy\n★-2014213605★"

[[96]]

[1] "eejongho: RT @na1m: 누리미디어, 빅데이터 추천기술 적용 학술정보서비스 오픈 https://t.co/...

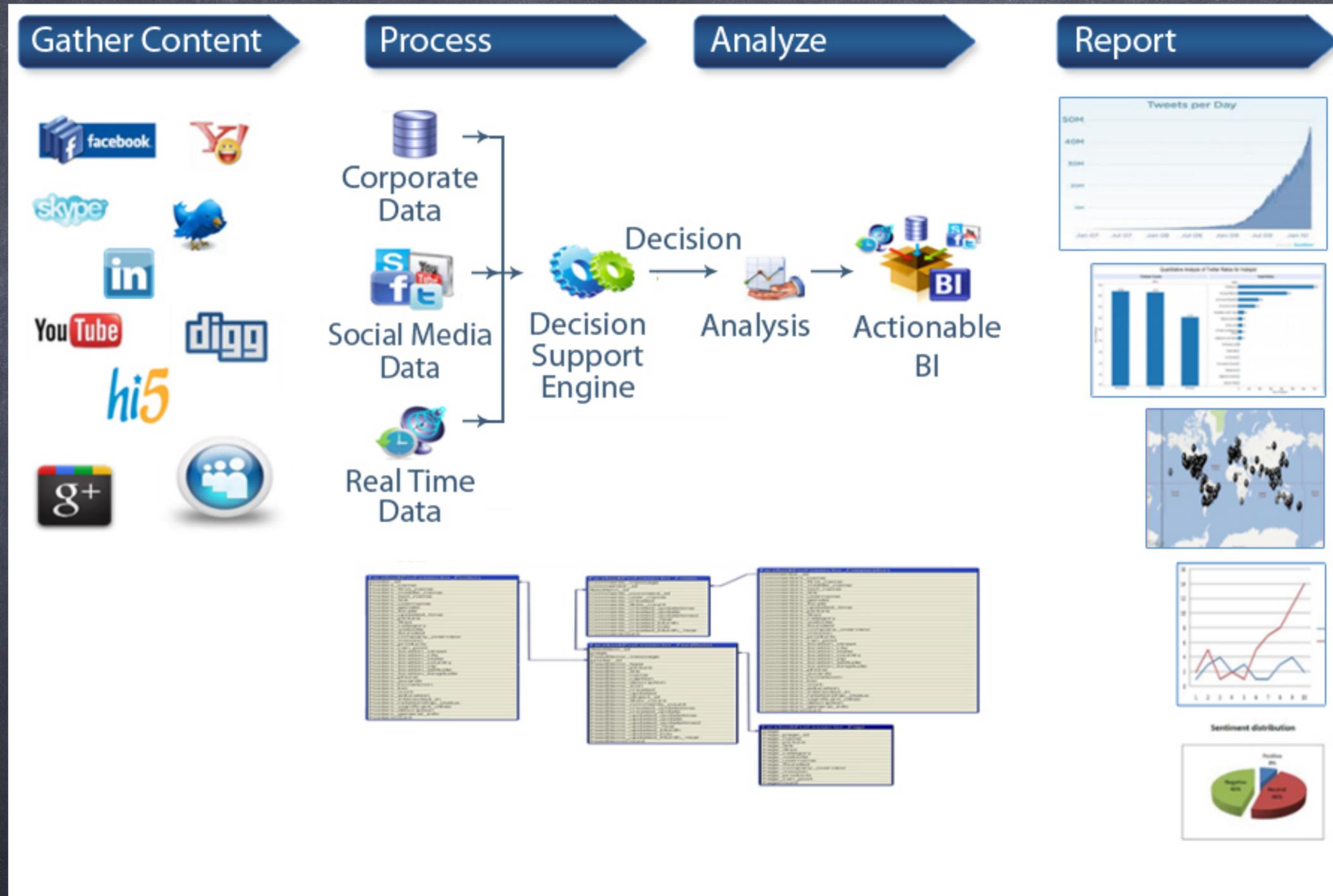
[[97]]

[1] "TheCatcherNews: 중국 정부가 의료 빅데이터 산업을 키우기 위해 나섰습니다 https://t.co/...
://t.co/X6bcdF88CE"

[[98]]

[1] "3RMarket0: [1편] ‘자원’ : 활용할 수 있는 빅데이터 발견하기\nhttps://t.co/Un7g3F...

Social Media Analytic



Social Media Analytic with Google Analytics

<http://analytics.google.com>

3단계로 사이트 트래픽 분석 시작

1 Google 애널리틱스 가입



모니터링하려는 사이트에 대한 기본 정보만 제공해 주시면 됩니다.

2 추적 코드 추가



귀하의 웹페이지에 설치해 두면 사이트에 방문자가 있을 때 Google에 알려주는 추적 코드가 제공됩니다.

3 잠재고객에 대해 알아보기



몇 시간 후부터 귀하의 사이트에 대한 데이터를 보실 수 있습니다.

Google Analytics

홈 보고서 맞춤설정 관리

wolfpack.hnu@gmail.com
한남대학교 통계학과 권세혁교수 - http://w...
전체 웹사이트 데이터

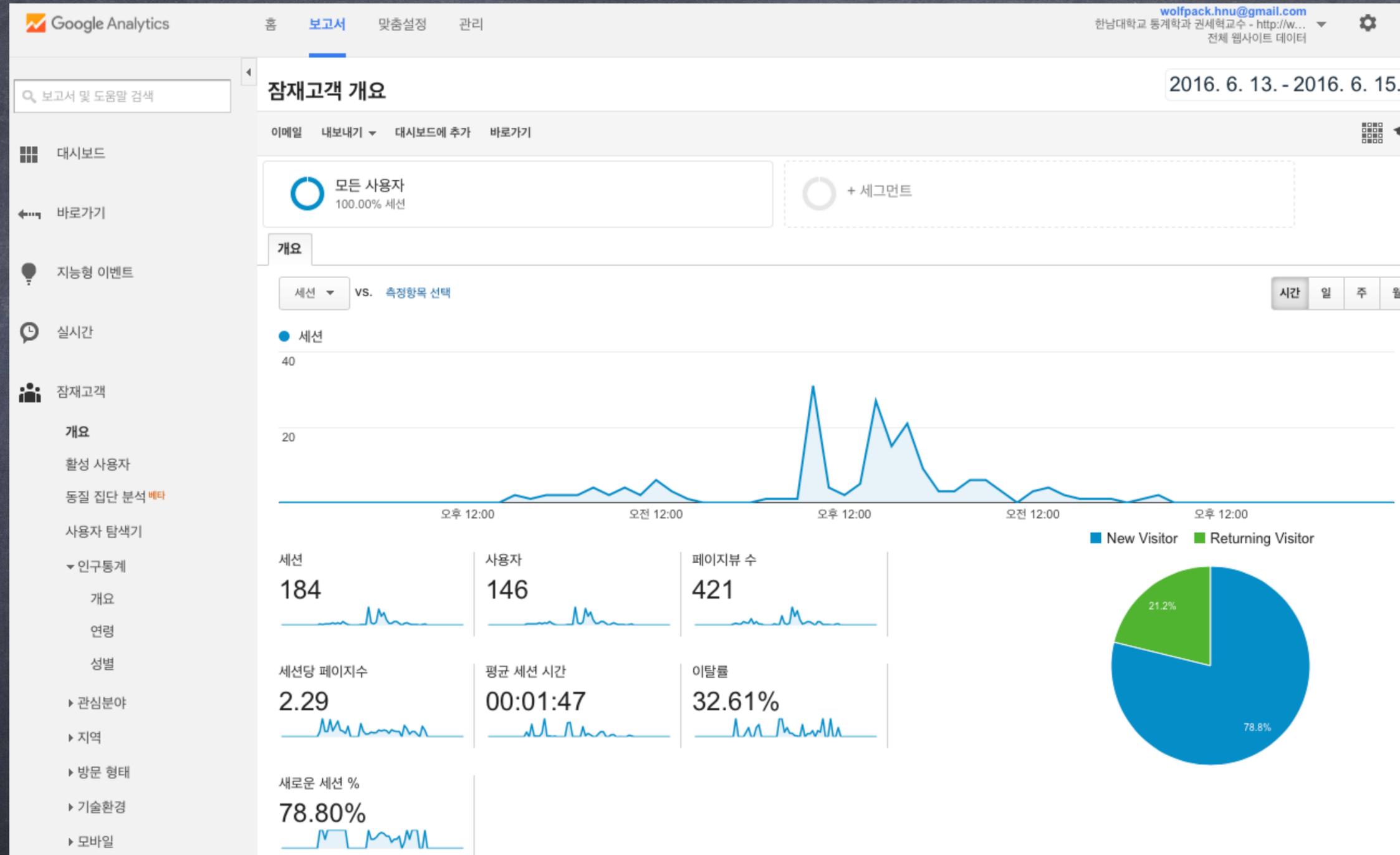
2016. 6. 15. - 2016. 6. 15.

↓ ↑

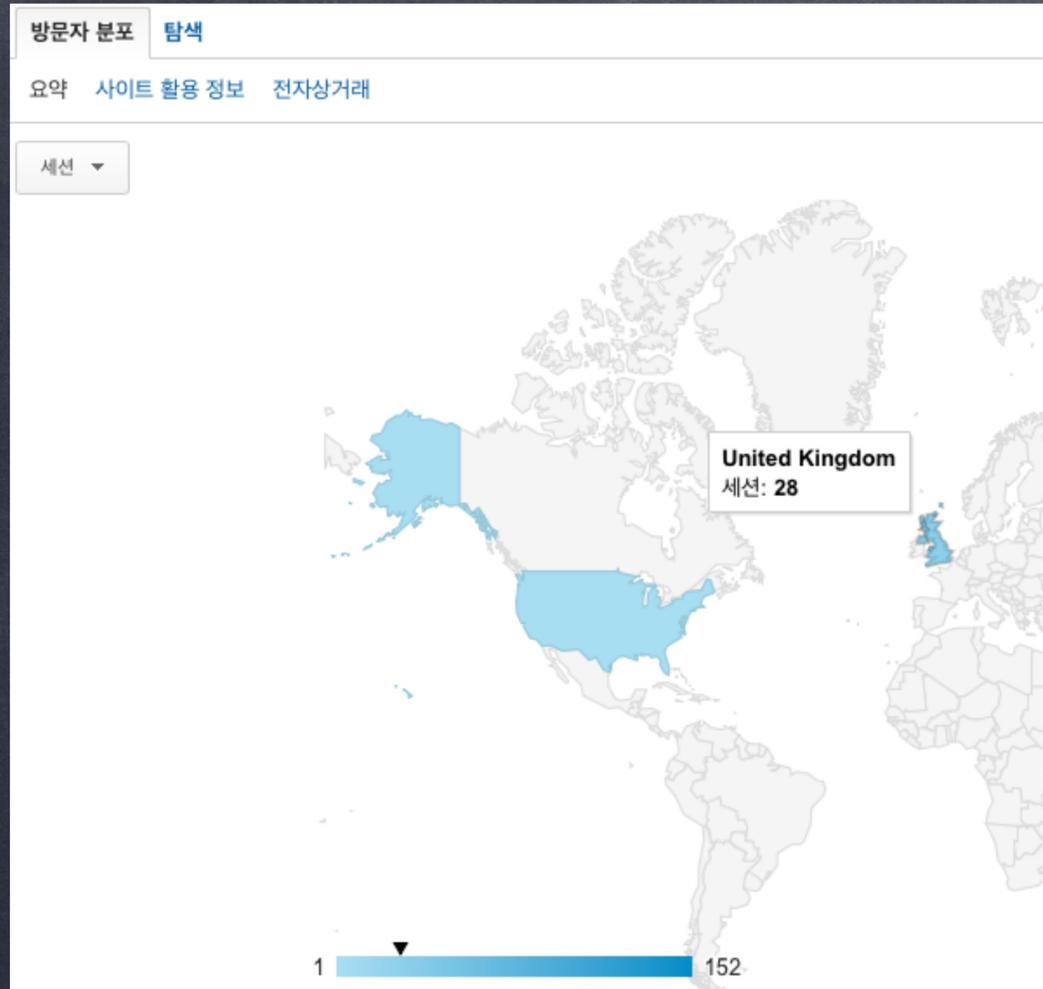
모드 표시 모두

세션	평균 세션 시간	이탈률	목표 전환율
☆ Wolfpack_HNU			
☆ 한남대학교 통계학과 권세혁교수 (UA-79234824-1)			

My webpages : Google analytics



My webpages : Google analytics (cont.)



지역	언어	위치	방문 형태	기술환경	모바일	개요	기기	맞춤	벤치마킹	사용자 흐름	획득			
											세션 ? ↓	새로운 세션 % ?	신규 방문자 ?	
<input type="checkbox"/>			신규 방문 vs. 재방문									57	63.16%	36
			계재빈도 및 방문빈도									전체 대비 비율 (%) : 30.81% (185)	평균 조화 : 78.38% (-19.42%)	전체 대비 비율 (%) : 24.83% (145)
<input type="checkbox"/>			참여도									28 (49.12%)	50.00%	14 (38.89%)
<input type="checkbox"/>												3 (5.26%)	33.33%	1 (2.78%)
<input type="checkbox"/>												2 (3.51%)	100.00%	2 (5.56%)
<input type="checkbox"/>												2 (3.51%)	50.00%	1 (2.78%)
<input type="checkbox"/>												2 (3.51%)	100.00%	2 (5.56%)
<input type="checkbox"/>												2 (3.51%)	50.00%	1 (2.78%)
<input type="checkbox"/>												2 (3.51%)	50.00%	1 (2.78%)
<input type="checkbox"/>												2 (3.51%)	50.00%	1 (2.78%)
<input type="checkbox"/>												2 (3.51%)	100.00%	2 (5.56%)
<input type="checkbox"/>												1 (1.75%)	100.00%	1 (2.78%)

Big data with R ggplot2()

```
install.packages("ggmap"); install.packages("ggplot2")
library(ggplot2); library(ggmap)

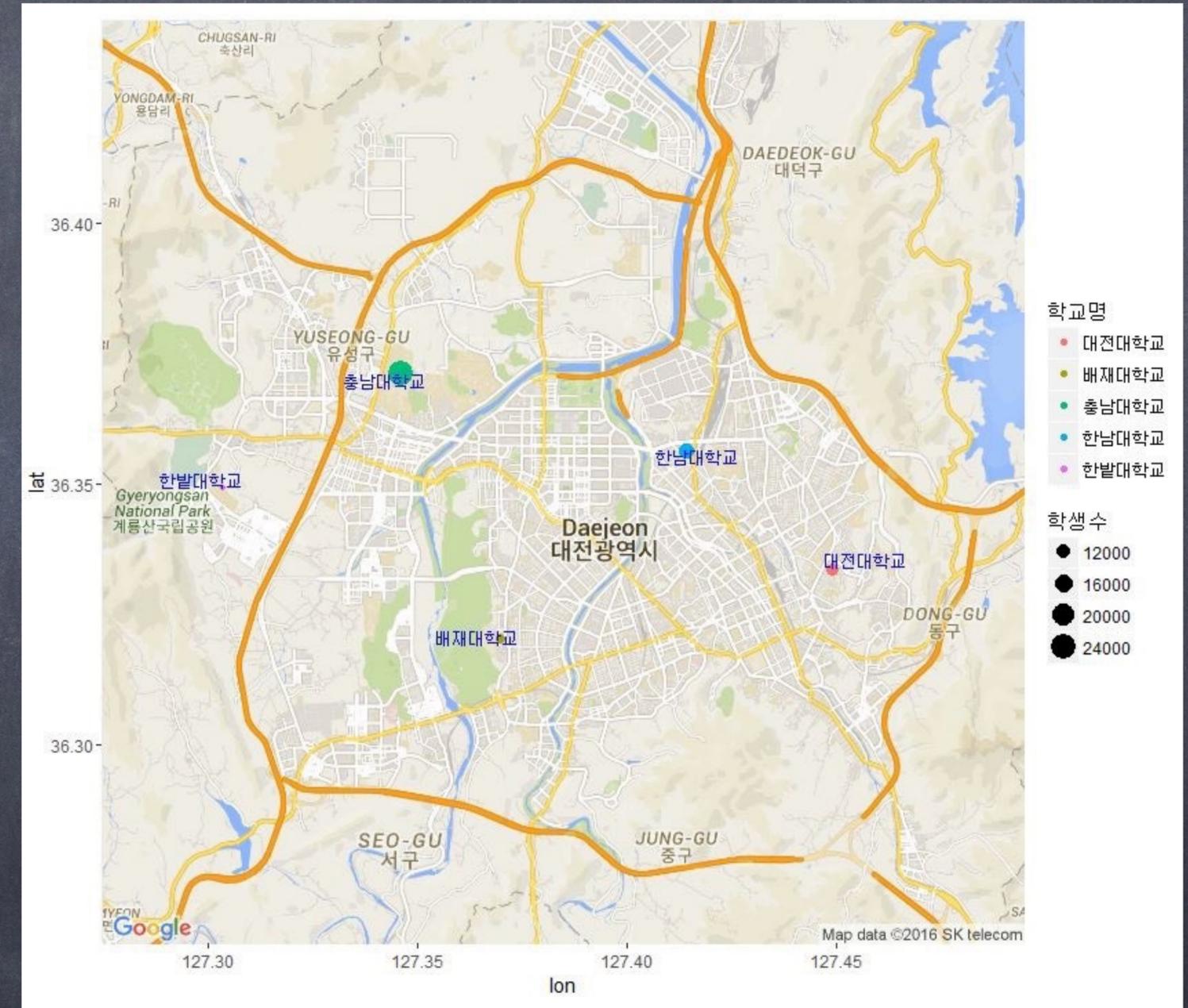
# maptypes=c("terrain", "terrain-background", "satellite",
# "roadmap", "hybrid", "toner", "watercolor", "terrain-labels", "terrain-lines",
# "toner-2010", "toner-2011", "toner-background", "toner-hybrid",
# "toner-labels", "toner-lines", "toner-lite")

uni_seoul <- read.csv("서울 11개 대학.csv", header=T)
seoul <- get_map("seoul", zoom=11, maptype = "roadmap")

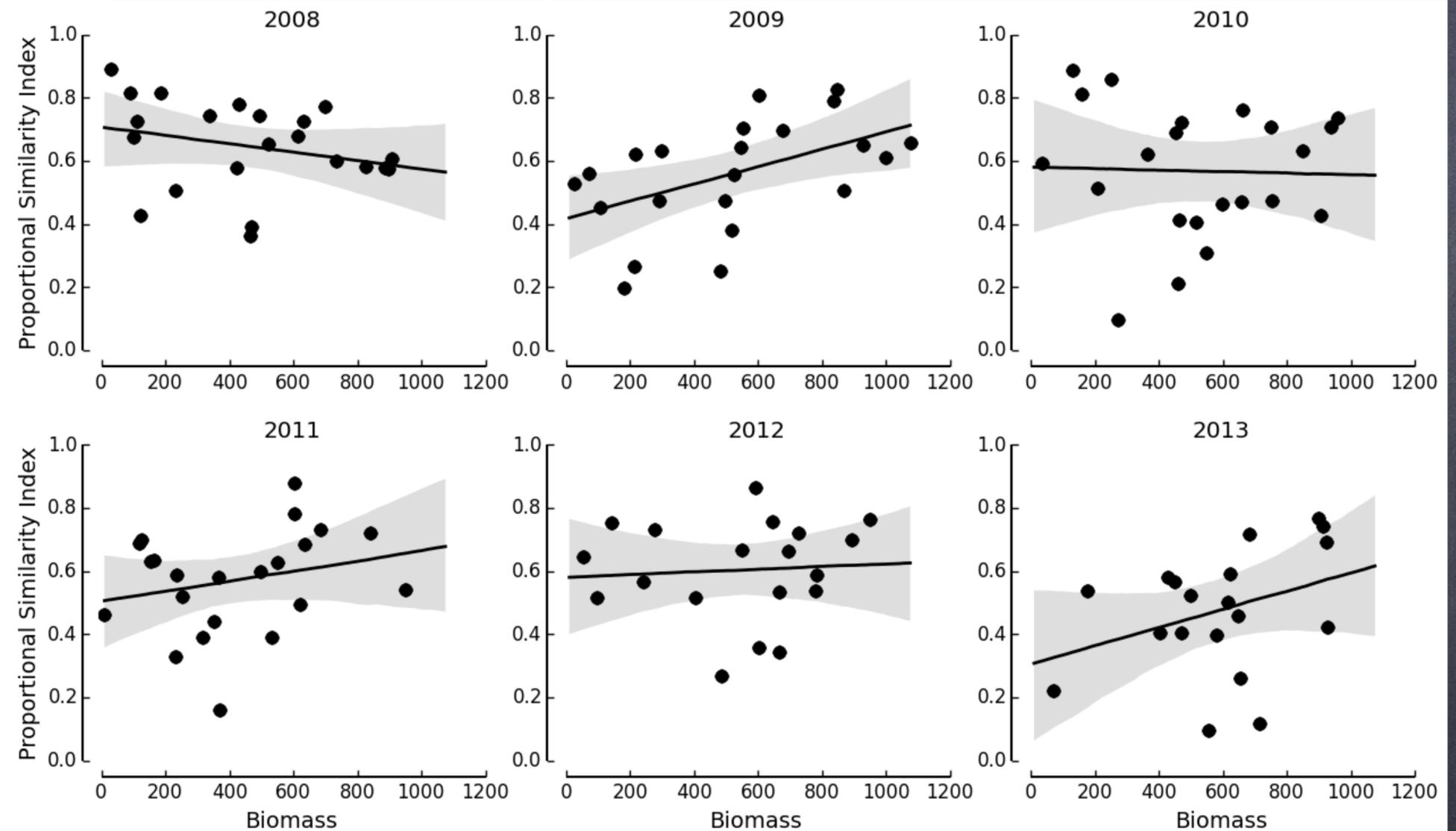
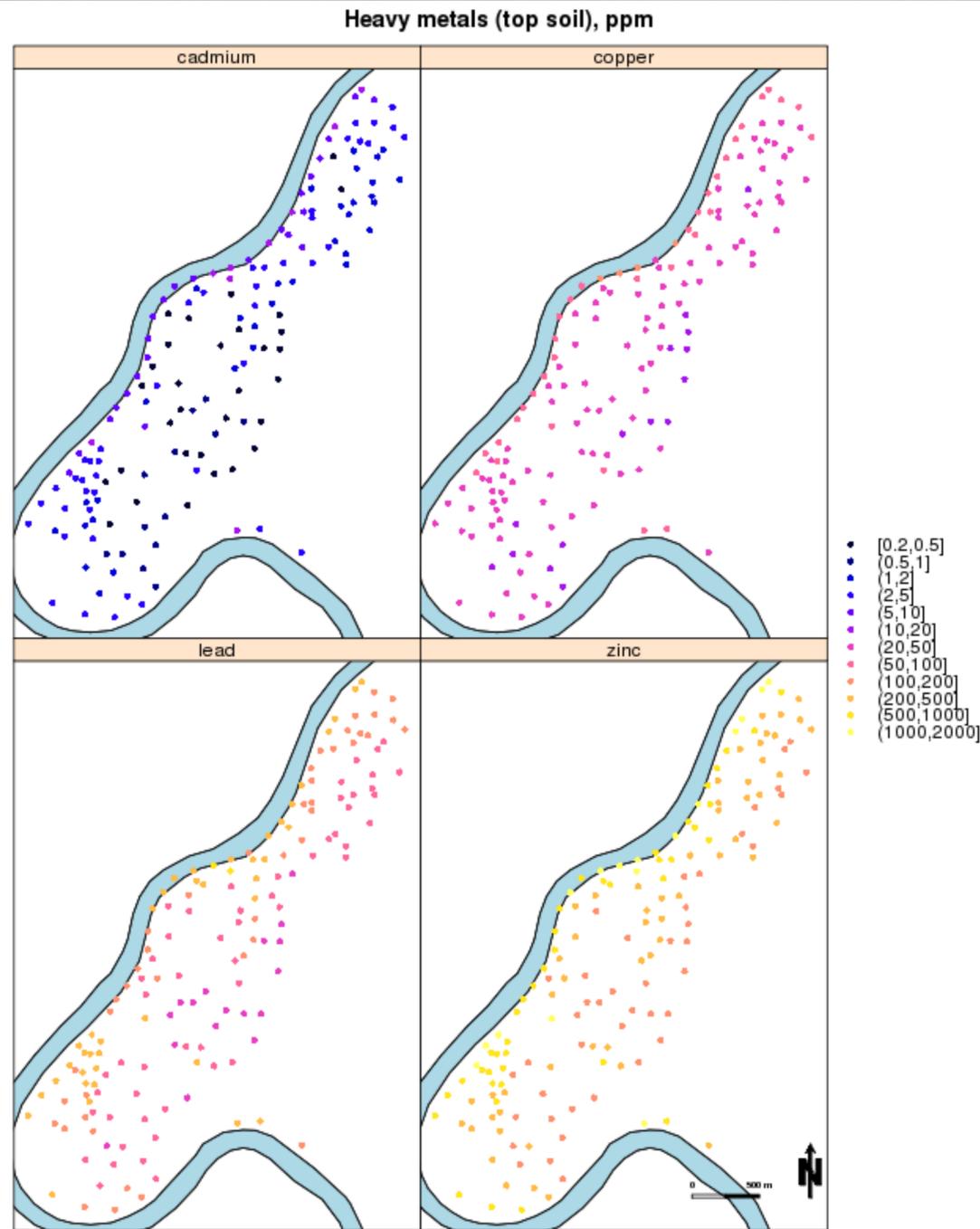
seoul_map <- ggmap(seoul)
seoul_map <- seoul_map + geom_jitter( data=uni_seoul, aes(x=경도, y=위도, size = 학생수,color=학교명)) +
scale_size(name="학생수")
seoul_map + geom_text(data=uni_seoul, aes(x = 경도, y = 위도, label=학교명),size=3,col="blue")
#=====
uni_daejeon <- read.csv("대전 5개 대학.csv", header=T)
daejeon <- get_map("daejeon", zoom=12, maptype = "roadmap")

daejeon_map <- ggmap(daejeon)
daejeon_map <- daejeon_map + geom_jitter( data=uni_daejeon, aes(x=경도, y=위도, size = 학생수,color=학교명)) +
scale_size(name="학생수")
```

ggplot2 결과



Big data with R lattice()



Still going...

R-bloggers

받는 사람: wolfpack_HNU

답장 받는 사람: R-bloggers

[R-bloggers] R, Yelp and the Search for Good Indian Food – An Open Course (and 6 more aRticles)

[R-bloggers] R, Yelp and the Search for Good Indian Food – An Open Course (and 6 more aRticles)

- [R, Yelp and the Search for Good Indian Food – An Open Course](#)
- [Using Microsoft R Server on a single machine for experiments with 600 million taxi rides.](#)
- [githubinstall: New R Package for Easy to Install R Packages on GitHub](#)
- [My knitr LaTeX template: manuscript and supplement interleaved in one source file](#)
- [R Hero saves Backup City with archivist and GitHub](#)
- [Le Monde puzzle \[#965\]](#)
- [R holds top ranking in KDnuggets software poll](#)

R, Yelp and the Search for Good Indian Food – An Open Course

Posted: 14 Jun 2016 03:09 PM PDT

(This article was first published on [DataCamp Blog](#), and kindly contributed to [R-bloggers](#))

New Free Course by Springboard and DataCamp

Are all Yelp restaurant reviews created equal? Should we place greater trust in reviews made by people who know the cuisine well? How about reviews of ethnically diverse cuisines? This free interactive tutorial will walk you through importing data, data manipulation, and data visualization. More importantly, it will teach you how to “cut the crap” from a large dataset.

[Play Course!](#)

This [free interactive tutorial](#) will walk you through importing data, data manipulation, and data visualization. More importantly, it will teach you how to “cut the crap” from a large dataset.

The screenshot displays a DataCamp course interface. On the left, a dark sidebar contains a notification: "Exercise Completed" with a "View" button. Below this, a message states: "Congratulations! You have finished the course and now know some good tools to manage data. You have also seen their work to solve an interesting problem. For more in-depth coverage on the concepts in this course by our Premium course..."

The main content area is split into two panes. The left pane shows R code for data visualization:

```
1 # The plotting package 'ggplot2' is
2 # Create a histogram of the avg_stars
3 hist(avg_review_indian$avg_stars)
4
5 # Create a histogram of the new_stars
6 hist(avg_review_indian$new_stars)
7
8 # Plot the distribution of changes to
9 hist(avg_review_indian$dif, main =
10 "Changes in Star Review", labs =
11 "Change")
12
13 # Plot the changes to per restaurant
14 ggplot(avg_review_indian, aes(x=1:nrow
15 (avg_review_indian), y=dif, fill=city
16 )) +
17   geom_bar(orientation="vertical", position
18   =position_dodge()) +
19   theme_minimal() + scale_fill_grey()
20   + labs("Business ID" = ylab("Change
21   In Star Review"))
22
```

The right pane shows a faceted histogram titled "Changes in Star Review". The y-axis is labeled "Change in Star Review" and ranges from -1 to 1. The x-axis is labeled "Business ID" and ranges from 0 to 20. The histogram is faceted by city, with a legend on the right listing: Annapolis, Chandler, Honolulu, Miami, Phoenix, San Antonio, and Tempe. The bars are colored in shades of grey.