

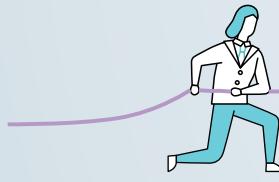
제7회

KOSSDA

데이터 페어



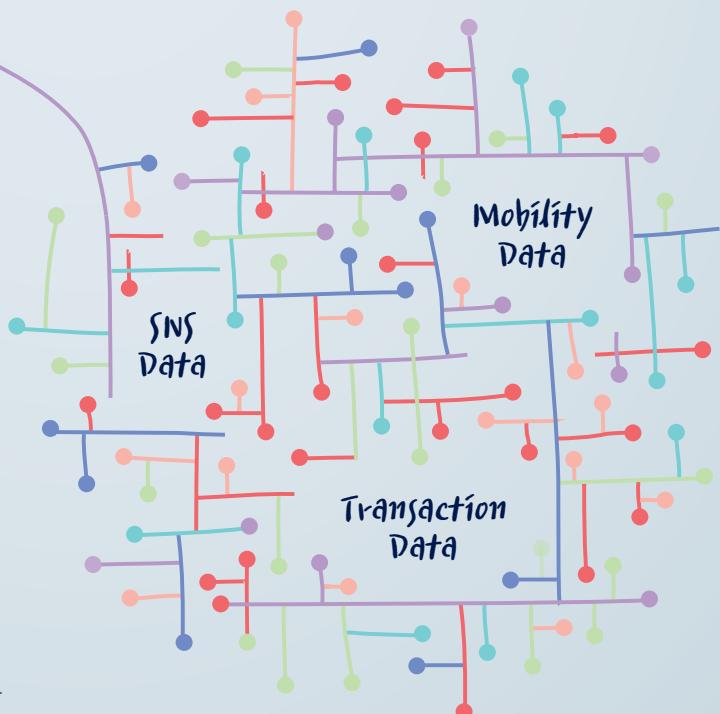
사회과학자가 빅데이터를 만나면 뭘 할까?



2019. 6. 27.

(목) 14:00-17:00

서울대학교
박물관(70동) 강당



KOSSDA 데이터 페어는 서울대 사회과학대학

'미래 기초학문분야 기반조성사업'의 지원을 받아 개최되는 행사입니다.

제7회

KOSSDA 데이터 페어

사회과학자가 빅데이터를 만나면

2019. 6. 27. (목) 14:00-17:00

서울대학교 박물관(70동) 강당

뭘 할까?

프로그램

I 인사말

이봉주 학장 (서울대 사회과학대학)

박수진 소장 (서울대 아시아연구소)

14:00-14:10

I 1부: 빅데이터 활용 방법론 강의

사회과학 분야 빅데이터 연구방법론

김기훈 교수 (서울대 사회학과, 사이람 대표)

14:10-14:50

휴식

14:50-15:00

I 2부: 빅데이터 활용 사례

모빌리티 데이터로 바라보는 사회

15:00-15:30

김정민 연구원 (카카오 모빌리티)

정치연구와 빅데이터

15:30-16:00

한규섭 교수, 노선헤 연구원 (서울대 언론정보학과)

문화권력에 대한 세 가지 데이터 분석

16:00-16:30

이원재 교수 (카이스트 문화기술대학원)

질의응답

16:30-17:00



제7회

KOSSDA

데이터 페어

목차

1부: 빅데이터 활용 방법론 강의

- 사회과학 분야 빅데이터 연구방법론 03
김기훈 교수 (서울대 사회학과, 사이람 대표)

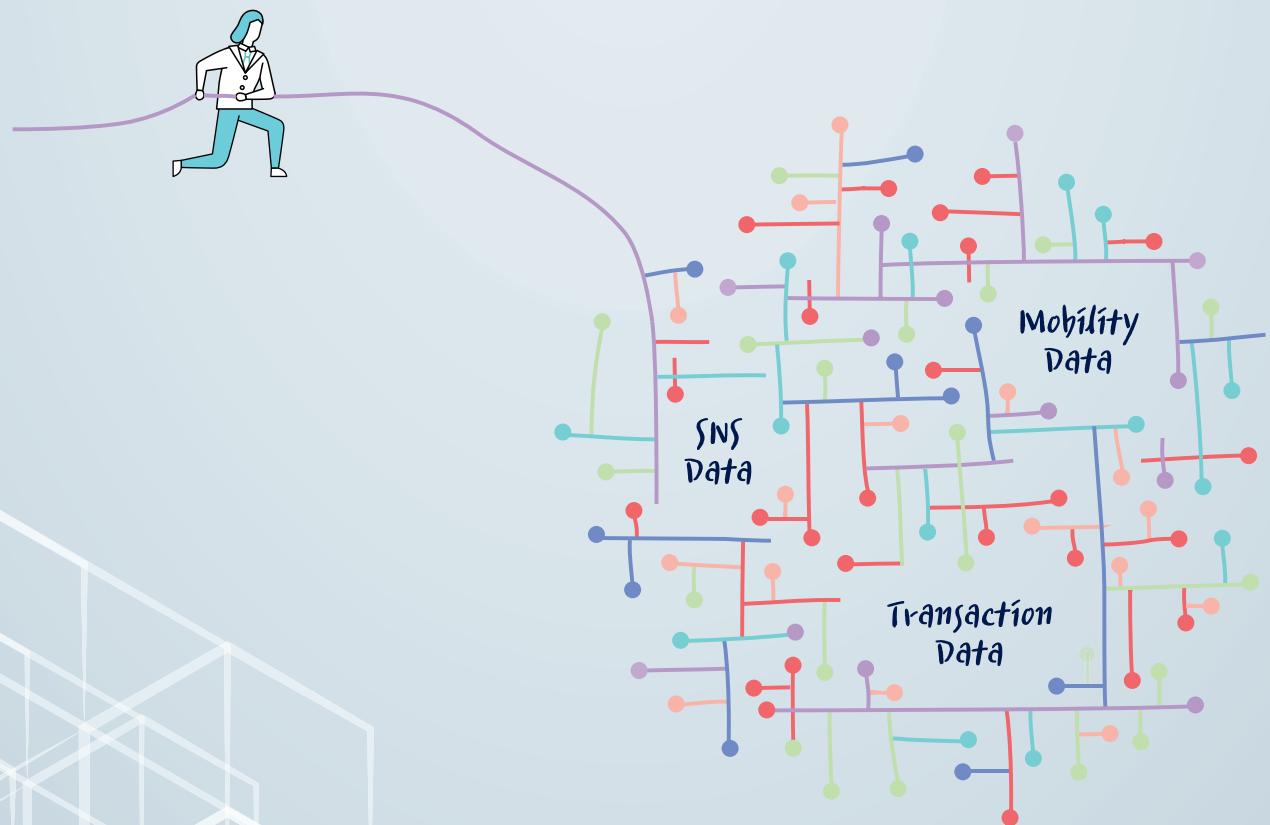
2부: 빅데이터 활용 사례

- 모빌리티 데이터로 바라보는 사회 39
김정민 연구원 (카카오 모빌리티)
정치연구와 빅데이터 57
한규섭 교수, 노선헤 연구원 (서울대 언론정보학과)
문화권력에 대한 세 가지 데이터 분석 81
이원재 교수 (카이스트 문화기술대학원)

▶ 1부: 빅데이터 활용 방법론 강의

사회과학 분야 빅데이터 연구방법론

김기훈 교수 (서울대 사회학과, 사이람 대표)





사회과학 분야 빅데이터 연구방법론

2019-6-27

김 기훈
(주)사이람 대표이사
서울대 사회학과



Table of Contents

- Big Data 의 본질
- Big Data 의 과학적 성격
- Big Data 의 Survey Error
- Big Data 의 사회과학적 의의
- Big Data 의 종류
- Big Data 의 수집
- Big Data 의 전처리
- Big Data 의 분석
- 사회과학 연구를 위한 Big Data 교육과정
- 사회과학 연구를 위한 Big Data 소프트웨어
- 맷음말

Big Data 의 본질

'Big Data' ?

'Big Data' as a misnomer(誤稱)?

Big Bang wasn't big.
There was no bang.

IT Features of "Big Data" : 3V

1. Volume
2. Variety
3. Velocity

'Big Bang'

3S
Sizable
Strong
Slow

essence/identity(種別의 本質) of 'Big Data' ?

'Big Animal'

→ "digital footprint(trace) data"

Elephant



Big Data

highly detailed “digital footprint/trace data”
automatically captured/generated using IT systems / sensors

Maciej Beresewicz et. Al.(2018), An overview of methods for treating selectivity in big data sources

Digital Footprint

"A digital footprint is the data trail left by interactions in a digital environment; including the use of TV, mobile phone, the World Wide Web, the internet and other connected devices and sensors.

Digital Footprints provide data on what has been performed in the digital environment (e.g. what you clicked on, searched for, liked, where you went, your location, your IP address, what you said, what was said about you)"



- Wikipedia -

"cyber shadow"
"electronic footprint"
"digital shadow"

Big Data 의 과학적 성격: 이론의 종말?

“데이터로 하여금 말하게 하라” -이론(Theory)의 종말?

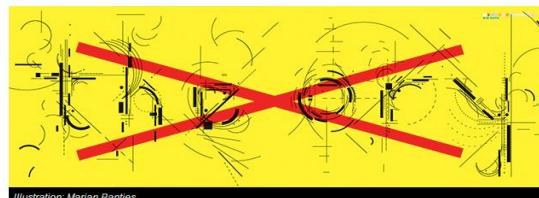


WIRED MAGAZINE: 16.07

SCIENCE / DISCOVERIES

The End of Theory: The Data Deluge Makes the Scientific Method Obsolete

By Chris Anderson 06.23.08



와이어드(Wired), 2008년 6월 23일자
커버스토리 '이론의 종말'
[크리스 앤더슨'](#)



“데이터로 하여금 말하게 하라” -이론(Theory)의 종말?

- 이제 세계는 엄청나게 많은 데이터와 응용 수학이 고려할 수 있는 다른 모든 도구들을 대체한다. 언어학에서부터 사회학에 이르기까지, 모든 인간 행동에 대한 이론은 버려라. 분류학, 존재론, 심리학도 잊어라. 어느 누가 사람들이 왜 무엇을 하는지 알까? 중요한 것은 사람들은 그것을(무엇인가를) 하고 있고, 우리는 사람들의 행위를 전례 없는 정확도로 추적하고 측정할 수 있다는 것이다. 충분한 양의 데이터와 함께라면, 숫자는 자기 생각을 말한다



“데이터로 하여금 말하게 하라” -이론(Theory)의 종말?

- 이렇게 큰 규모의 컴퓨터 사용을 배우는 것은 매우 도전적일지 모른다. 그러나 기회는 대단할 것이다. 엄청난 양의 데이터, 그리고 이러한 숫자를 다룰 수 있는 통계적 도구는 세계를 이해하기 위한 완전히 새로운 방식을 제공한다. 상관성(Correlation)이 인과성(Causation)을 대체하고 과학은 일관적인 모델, 통합된 이론, 혹은 설명 기제 없이도 진보할 것이다.

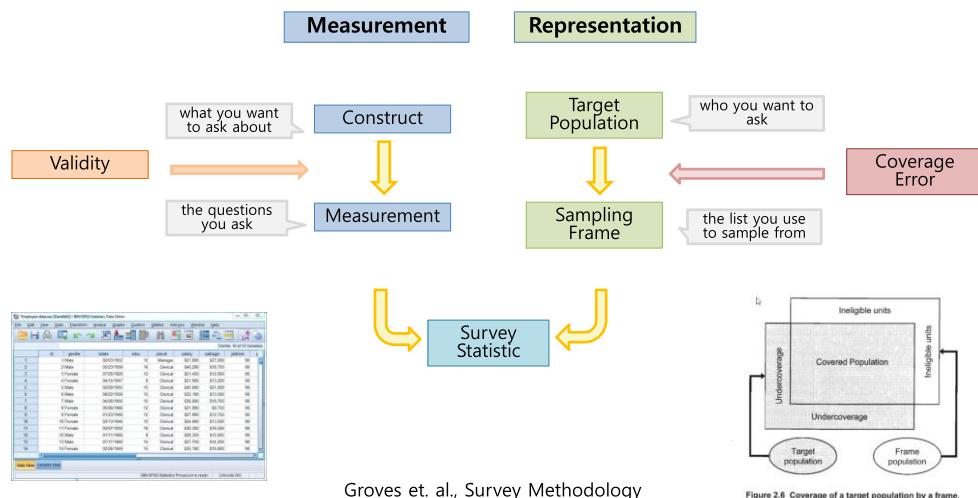
Big Data vs Survey Error

Analytic features of Big Data

- not designed for (statistical) analysis ("organic data", "administrative data")
- non-probabilistic character
- mediation of digital platform
- unstructured / highly detailed

Maciej Beresewicz et. Al.(2018), An overview of methods for treating selectivity in big data sources

Survey Error



Groves et. al., Survey Methodology

Figure 2.6 Coverage of a target population by a frame.

Big Data 의 Survey Error

	cons	pros
representation /coverage of population	non-probabilistic sample data → coverage error	coverage of 'entire' population
validity /measurement	mediation of digital platform → measurement error	unobtrusive / nonreactive (cf. mediation of questionnaire) "설문조사는 그 자체가 하나의 외부적 요인으로서, 우리가 묘사하고자 하는 사회적 상황을 침범 (intrude)함으로써, 어떤 태도를 측정하고 동시에 형성"
unit & variable of interest	<ul style="list-style-type: none"> unit = individual → unit error variable = attribute → proxy variable <p>→ "unstructured"</p>	<ul style="list-style-type: none"> unit = event (→ social action/interaction) variable = (structural) position <p>→ "detailed"</p>

Big Data 의 Coverage :

연구의 Target population(A)과 해당 Big Data 를 생성한 특정 digital platform(B)과의 관계

관계의 종류	A의 예 (B: Twitter)	Big data 분석의 의미	비고
$A = B$	“Twitter”	complete enumeration survey	단, 해당 플랫폼의 데이터 수집 방법에 의존적 기술통계, 기계학습에 의거한 예측모델링. (한정된 경계 내의)
$A \supset B$	“online social network”	case study	소규모(개인, 소집단) → 초대규모(전체 디지털 플랫폼) (참여)관찰 → (객관적)기록(log) 수작업 → 자동화 심층(in-depth) → 상세(detail) 질적 분석 → 양적 분석 기술적 → 기술적(+탐색적/확증적)
$A ? B$	“한국 국민”	non-probabilistic (sample) survey	with proper method of calibration /adjustment (using auxiliary variable) → “big data survey”(?)

Big Data 의 사회과학적 의의



Cameron Marlow

I am a research scientist and "in-house sociologist" at Facebook. My research focuses on various aspects of online communities including the diffusion of information across online social networks, access to information and social capital, and the incentives that impact social media production.

I received my Ph.D. at the MIT Media Laboratory on the topic of media contagion and weblogs. You might also be familiar with my project Blogdex which tracks (used to track?) diffusion in the weblog community. The results of this work and the MIT Weblog Survey can be found in [my Ph.D. Thesis](#).

After finishing my Ph.D. I joined Yahoo! Research Berkeley where I ran the Social Motives group studying the social incentives in emerging applications such as Flickr, del.icio.us, and Last.fm. I joined Yahoo! Research and continued my work on various topics related to social search, community evolution, and network topology.

In addition to my research, I currently manage the Data Science team at Facebook, and am engaged in trying to understand how the data platforms of the future will change the roles of researchers in the enterprise and academia.

email: cam@alum.mit.edu
 aim: cameronfactor
[flickr](#)
[facebook](#)
[linkedin](#)

Recent Papers

- [Structural diversity in social contagion](#)
- [The Role of Social Networks in Information Diffusion](#)
- [Why Most Facebook Users Get More Than They Give: The Effect of Facebook "Power Users" on Everybody Else](#)
- [A 61-million-person experiment in social influence and political mobilization](#)
- [Social capital on Facebook: Differentiating uses and users](#)

Facebook의 Data Science Team은
사회학자가 진두지휘 하고 있음.

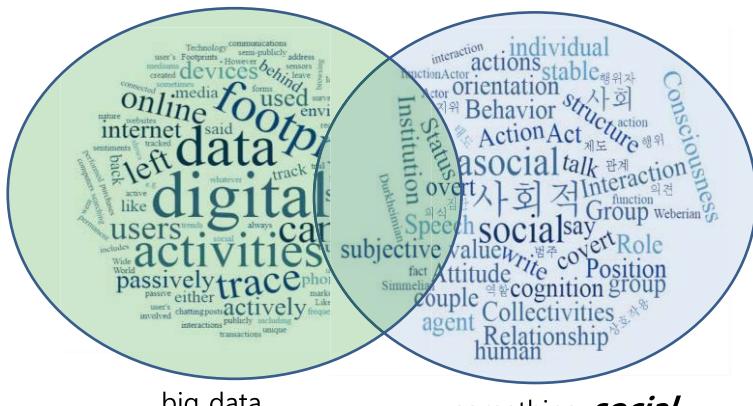


이제 우리는 과거에는 절대 불가능했던 정도의 초정밀현미경으로 인간의 사회적 행위를 들여다 볼 수 있게 되었다. 나아가 수백만 명을 상대로 실험을 하는 것조차 가능하다.

Facebook이 풀어야만 하는 최대의 과제는 사회과학이 풀어야 하는 과제와 똑같은 것이다.

A spectre is haunting social sciences – the spectre of big data.

social big data



big data
digital footprint data

something **social**

- Weberian social action
- Simmelian sociation
- Durkheimian social facts

소셜 빅데이터의 사회과학 방법론적 성격

	전통적 사회조사	소셜 빅데이터
Social	Proxy social (attributes of actors)	Non-proxy social <ul style="list-style-type: none"> social action/interaction speech act
Obtrusiveness	Obtrusive (questionnaire/interview)	Non-obtrusive (observation)
Sampling	Sampling	Non-sampling (complete enumeration)

the “Social”

Max Weber
(1864-1920)



Social Action

Georg Simmel
(1858-1918)

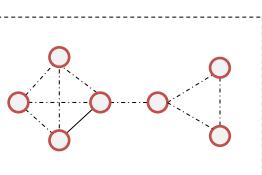


Sociation

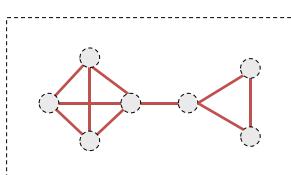
Emile Durkheim
(1858-1917)



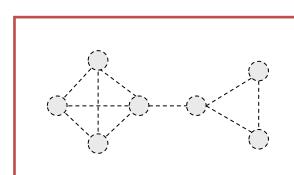
Social Fact



consideration/orientation
meaning



Geometry
position



sui generis,
emergent properties

Rediscovery of hidden dimension in Big Data : "social action" as data

"The devil is in the detail"

Big Data = detailed transaction data left in a digital environment

"an action is '**social**' if the acting individual takes account of the behavior of others and is thereby oriented in its course".

– Max Weber

Max Weber
(1864-1920)



Social action

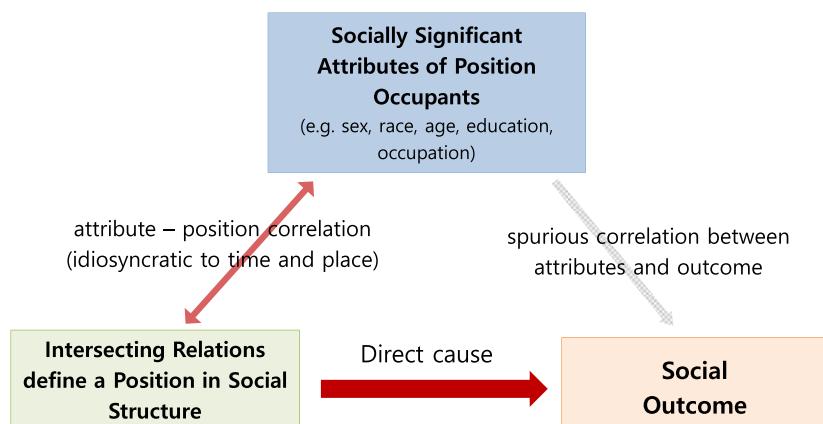
Georg Simmel
(1858-1918)



Sociation

'sleeping beauties' in detail fields of big data
: social action, *Verstehen* (Weber)
: sociation (Simmel)

Rediscovery of hidden dimension in Big Data : "social action" as data



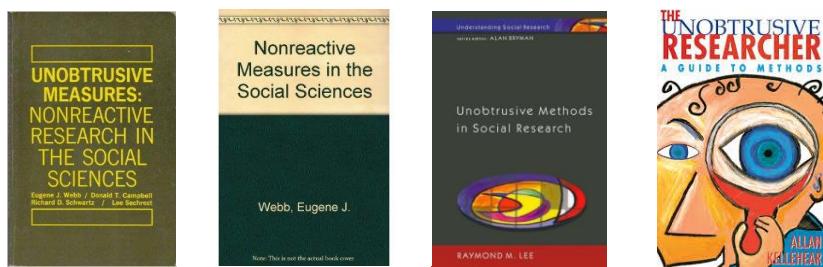
Ronald S. Burt. (1991). Structure : Reference Manual, p. 13

unobtrusive/nonreactive data (무반응/불개입 데이터)

"오늘날 사회과학 연구조사는 인터뷰 혹은 설문조사를 기초로 한다. 우리는 이런 오류가 있을 수 있는 단일한 방법에 지나치게 의존하는 것을 어렵게 생각한다.

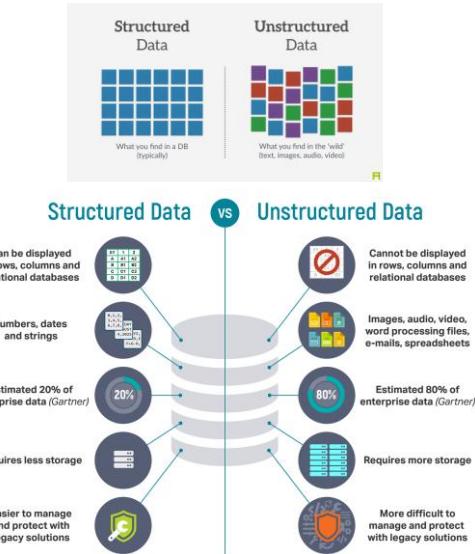
인터뷰와 설문조사는 그 자체가 하나의 외부적 요인으로서, 우리가 묘사하고자 하는 사회적 상황을 침범(intrude)함으로써, 어떤 태도를 측정하고 동시에 형성한다. 또한 그러한 방법들은 응답에 협조적이고, 접근 가능한 사람에게만 사용될 수 있을 뿐만 아니라, 조사 토픽과 관계없는 개별적 차원에 의해 형성된 응답일 수 있다는 한계점을 가진다.

그러나, 가장 주된 반대(objection)이유는 이런 방법들이 단독으로 사용되어 서는 안 된다는 것이다."



Big Data의 종류

Structured data vs. Unstructured data



Master Data vs. Transaction Data

Transaction Data vs. Master Data						
Transaction Data						
Customer	Date	Product	Code	Price	Quantity	Location
Stefan Kraus	1/2/2017	Scarpa Telemark Ski Boot	SC1279	€250	1	St. Moritz, CH
Donna Burbank	1/5/2017	Scarpa Telemark Ski Boot	SCU1289	\$150	1	Boulder, CO
Stefan Kraus	1/2/2017	North Face Down Jacket	NF8392	€450	1	Zurich, CH
Stefan Kraus	1/2/2017	Garmin Sports Watch	GM29384	€200	2	Zurich, CH
Wendy Hu	3/4/2017	Prana Yoga Pant	PN82734	\$51	5	New York, NY
Joe Smith	4/1/2017	Garmin Sports Watch	GM29384	\$150	1	Albany, NY

Master data represents the business objects that contain the most valuable, agreed upon information shared across an organization.^[1] It can cover relatively static reference data, transactional, unstructured, analytical, hierarchical and metadata.^[2] It is the primary focus of the information technology (IT) discipline of master data management (MDM).

Master data is usually non-transactional in nature, but in some cases gray areas exist where transactional processes and operations may be considered master data by an organization. For example, master data may contain information about customers, products, employees, materials, suppliers, and vendors. Though rare, if that information is only contained within transactional data such as orders and receipts and is not housed separately, it may be considered master data.^[3]

Transaction data is data describing an event (the change as a result of a transaction) and is usually described with verbs. Transaction data always has a time dimension, a numerical value and refers to one or more objects (i.e. the reference data).

Typical transactions are:

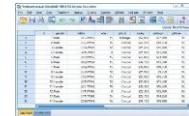
- Financial: orders, invoices, payments
- Work: plans, activity records
- Logistics: deliveries, storage records, travel records, etc.

Big Data 의 3 종류 (3T)

빅데이터 중 팔할은 텍스트 데이터 or 트랜잭션 데이터
(텍스트 데이터 반, 트랜잭션 데이터 반)

1. Traditional Data

(개체-속성 데이터)



2. Text Data



3. Transaction Data

ID	NAME	AGE	ADDRESS	BIRTHDAY	EDUCATION	JOBCAT	SALBEGIN	JOBTIME	P
1	John Doe	35	New York, USA	1985-01-01	High School	XX	\$100000.00	Full-time	98
2	Jane Doe	32	New York, USA	1988-12-31	College	XX	\$100000.00	Full-time	98
3	Mike Johnson	40	Los Angeles, USA	1985-01-01	High School	XX	\$100000.00	Full-time	98
4	Sarah Johnson	38	Los Angeles, USA	1988-12-31	College	XX	\$100000.00	Full-time	98
5	David Williams	30	Chicago, USA	1990-01-01	College	XX	\$100000.00	Full-time	98
6	Alice Williams	28	Chicago, USA	1992-12-31	College	XX	\$100000.00	Full-time	98
7	Bob Smith	35	Seattle, USA	1985-01-01	High School	XX	\$100000.00	Full-time	98
8	Mary Smith	33	Seattle, USA	1988-12-31	College	XX	\$100000.00	Full-time	98
9	Tom Jones	42	Atlanta, USA	1985-01-01	College	XX	\$100000.00	Full-time	98
10	Linda Jones	40	Atlanta, USA	1988-12-31	College	XX	\$100000.00	Full-time	98
11	Steve Parker	38	Boston, USA	1985-01-01	College	XX	\$100000.00	Full-time	98
12	Karen Parker	36	Boston, USA	1988-12-31	College	XX	\$100000.00	Full-time	98
13	Paul Parker	34	Boston, USA	1990-01-01	College	XX	\$100000.00	Full-time	98
14	Emily Parker	32	Boston, USA	1992-12-31	College	XX	\$100000.00	Full-time	98

(1) Traditional Data: 개체-속성 데이터 (master data)

Employee data.sav [DataSet1] - IBM SPSS Statistics Data Editor									
	id	gender	bdate	educ	jobcat	salary	salbegin	joftime	p
1	1	Male	02/03/1952	15	Manager	\$57,000	\$27,000	98	
2	2	Male	05/23/1958	16	Clerical	\$40,200	\$18,750	98	
3	3	Female	07/26/1929	12	Clerical	\$21,450	\$12,000	98	
4	4	Female	04/15/1947	8	Clerical	\$21,900	\$13,200	98	
5	5	Male	02/09/1955	15	Clerical	\$45,000	\$21,000	98	
6	6	Male	08/22/1958	15	Clerical	\$32,100	\$13,500	98	
7	7	Male	04/26/1956	15	Clerical	\$36,000	\$18,750	98	
8	8	Female	05/06/1966	12	Clerical	\$21,900	\$9,750	98	
9	9	Female	01/23/1946	15	Clerical	\$27,900	\$12,750	98	
10	10	Female	02/13/1946	12	Clerical	\$24,000	\$13,500	98	
11	11	Female	02/07/1950	16	Clerical	\$30,300	\$16,500	98	
12	12	Male	01/11/1966	8	Clerical	\$28,350	\$12,000	98	
13	13	Male	07/17/1960	15	Clerical	\$27,750	\$14,250	98	
14	14	Female	02/26/1949	15	Clerical	\$35,100	\$16,800	98	

(2) Text 데이터

주제
 본 보고서는 디자이너, 구현, 이미지, 동영상, 사용도, 버포, 사용, 풍부시 소셜미디어를 연동하고, 소셜미디어 통합 계정정보와 콘텐츠 발급, 사용, 공유, 풍부한 콘텐츠의 사용정보 수집을 통한 관리, 사용경로! 리 시스템 서비스(100)와 소셜미디어 시스템 서비스(200)를 연동하여 기반 소셜미디어 활용률 확장! 디자인 기반 활용률 확장! 디자인 기반 소셜미디어 활용률 확장! 디자인 기반 활용률 확장! 디자인 기반 활용률 확장! 디자인 기반 활용률 확장! 디자인 기반 활용률 확장!

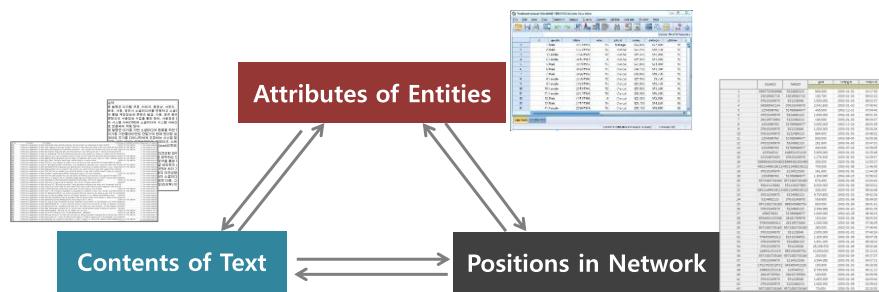
Created Time	Writer	APP	HashTags	RT Count	Full Text	Type
2018-12-07 19:00: sophiakaito	Twitter for i None	0	@hi_ром974	71	졌다... 언니는 추워서... 겨울 스카프 이 날씨에 맞아요....	Reply
2018-12-07 19:00: seojohn_nctzen	Twitter for i NCT	2302	RT @limehane	181207 SME 에스 엠 터미널트 사과공연 축운 날씨에도 NCNT 불사랑을 즐기자	RT	
2018-12-07 19:00: yuini19_tao	Twitter for i 자연Dream	4040	RT @official_zone	FNS 기요제에 참가해주시면서 감사합니다며 정말 많이 열었던 무대였습니다.	RT	
2018-12-07 19:00: HS_CB	Twitter for i 농형식.park1	50	RT @sik_shine	181207 뮤지컬 <엘리자벳> 최근 날씨가 너무 추워지니까 솔직히 유금님 유금님	RT	
2018-12-07 19:00: wej053161	Twitter for i None	495	RT @Aseyeon077	광주티비에서 이세영 제작 첫 시위공연 시장을 거쳐 광주고속버스터미널	RT	
2018-12-07 19:00: studyr_feather	Twitter for i None	7311	RT @BTS_twt	안녕하세요 오늘 날씨가 다시 추워진듯 해요 다음 찾으셨던가고 다니	RT	
2018-12-07 19:00: wjlover_	Twitter for i None	528	RT @Babybird_Jo	영하의 추운 날씨에 오전부터 일렬한 우리들을 위해 고기 반찬 기득가득 맛있	RT	
2018-12-07 19:00: mongmong9	Twitter for i None	9464	RT @seulgiyou	우연히 산 증고책, 세월이 흘러 태우리가 노랗게 바랜 첫 장에 23년 전 누군가의	RT	
2018-12-07 19:00: lsggg	Twitter for i None	242	RT @hyang0404	<파릇한 국가인재원> 내 눈에서 눈물나게 한 사람을, 취임 100일이라고 직원	RT	
2018-12-07 19:00: bamjswz	Twitter for i None	50	RT @daystar	990527 친환경 추운 날씨에 이쁜 아침부터 수고했어 체리딸기왕님 https://t.co/	RT	
2018-12-07 19:00: Mabi_Seanan	twittbot.net	0	내일 날씨는 어때?		Tweet	
2018-12-07 19:00: jumjye_0621	Twitter for i None	9464	RT @seulgiyou	우연히 산 증고책, 세월이 흘러 태우리가 노랗게 바랜 첫 장에 23년 전 누군가의	RT	
2018-12-07 19:00: drachane	Twitter for i None	528	RT @Babybird_Jo	영하의 추운 날씨에 오전부터 일렬한 우리들을 위해 고기 반찬 기득가득 맛있	RT	
2018-12-07 19:00: red_bsrglasses	Twitter for i 97페.낙태:29	29	RT @safe_abortion	<한법제판소> 낙태의 위험 판장을 축구하는 원법제판소 앞 19시위 #97페	RT	
2018-12-07 19:00: movement_0529	Twitter for i None	50	RT @daystar	990527 친환경 추운 날씨에 이쁜 아침부터 수고했어 체리딸기왕님 https://t.co/	RT	
2018-12-07 19:00: Bboing_39	Twitter for i None	0	진짜 오늘은 그냥 나가면 죽을 수도 있겠나를 생각하게 하는 날씨		Tweet	
2018-12-07 19:00: tahongthuongg	Twitter for i 재현JAEHY	445	RT @mydarlinghyun	181113 날씨가 많이 흡네요 웃 떠오르게 입으로세요 #재현 #JAEHYUN #NRT	RT	

(3) Transaction 데이터

순번	보험종목	회사	사고발생일자	지급일자	병원이름	사유	피보험자주민번호	피보험자	계약자주민번호	계약자	수익자	설계사
1	생명보험	삼성생명	1998-12-30	1998-12-31	서울대병원	OO	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아니	나아니	최고임
2	생명보험	현대해상	1998-12-31	1999-01-01	연세대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아니	나아니	최고임
3	자동차보험	MG	1999-01-01	1999-01-02	순천향대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아니	나아니	최고임
4	상해보험	메리츠	1999-01-02	1999-01-02	서울대병원	DD	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아니	나아니	최고임
5	생명보험	삼성생명	2000-10-20	2000-10-20	연세대병원	AA	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아니	최고임
6	생명보험	현대해상	2000-10-20	2000-10-20	가천의대병원	BB	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
7	자동차보험	MG	2000-10-20	2000-10-20	서울대병원	CC	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
8	생명보험	삼성생명	2003-04-05	2003-04-05	연세대병원	EE	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
9	생명보험	현대해상	2003-04-05	2003-04-05	인하대병원	FF	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
10	자동차보험	MG	2003-04-05	2003-04-05	차병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
11	상해보험	메리츠	2003-04-05	2003-04-05	연세대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	나아니	최고임
12	생명보험	삼성생명	2010-02-03	2010-02-03	순천향대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아니	최고임
13	생명보험	삼성생명	2019-02-19	2019-02-19	서울대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
14	상해보험	메리츠	2019-02-19	2019-02-19	연세대병원	XX	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임



	Traditional	Text	Transaction
Object	Entity	Contents	Relationship
Focus	Association among Attributes of Entities	Meanings contained in the Text	Structural Positions in Network
Analytic Methodology	Multivariate Statistics Machine Learning	NLP Text Mining	Network Analysis Graph Mining



Big Data 의 수집

빅데이터의 일반적 획득 방법

데이터 획득 방법	장점	단점
데이터 제공로부터 Data File 획득 (TXT, CSV 또는 엑셀 파일 형태)	직접 수집할 필요가 없음	<ul style="list-style-type: none"> 정보가 변경되었을 때 실시간으로 전달할 방법이 없음 시간 비용 소모가 큼
Database 접근권한을 얻어 데이터 수집	정형화된 데이터 수집이 용이 함	보안의 위험성(현관문 비밀번호를 알려준 것과 동일)
데이터 제공 업체의 API를 통한 데이터 수집	보안과 실시간 전달의 문제점을 해결	데이터 수집 제한이 있음
Web Crawling 통한 데이터 수집	거의 대부분의 웹 상의 데이터를 수집 가능	<ul style="list-style-type: none"> 웹 상의 소스코드를 분석하고 parsing 하는 것이 어려움 저작권 위반의 위험

WEB crawling

news.naver.com

HTML source

문자열 Data

HTML 페이지를 가져옴
 python
 Parsing(특정 문자 추출)을 통해 원하는 문자 정보를 얻는 것

문자열 데이터를 변수에 넣고

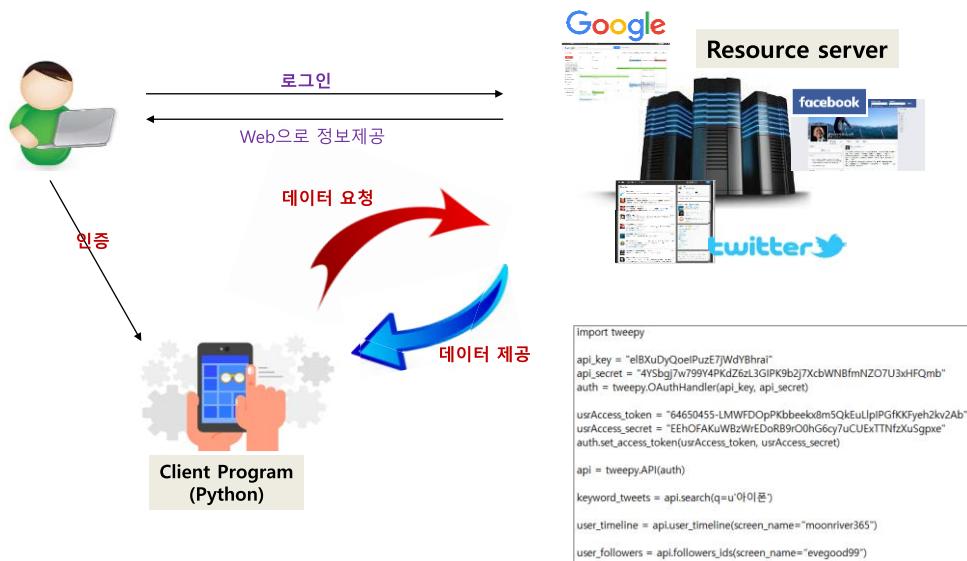
```

urlOpen = urllib2.urlopen("https://www.clien.net/service/board/sold")
html = urlOpen.read()

findex = 0
titleFirst = '<span data-role="list-title-text title="'
titleLast = ">"'

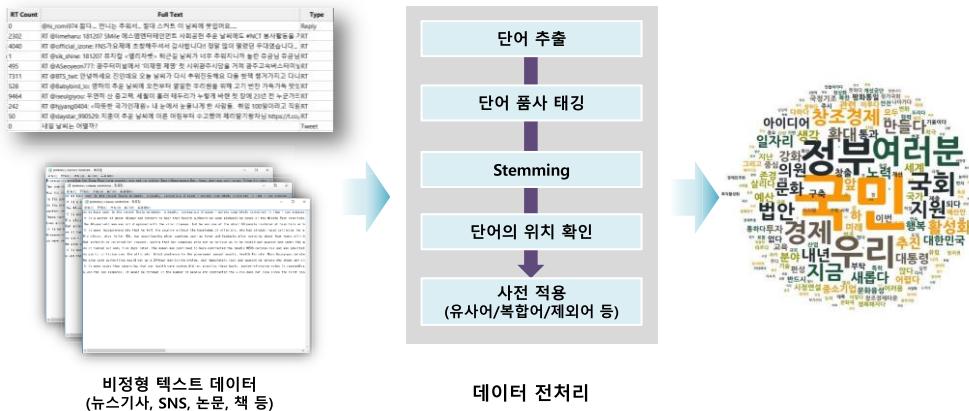
while True:
    try:
        titleFirstIndex = html.index(titleFirst, findex)+len(titleFirst)
    except:
        break
    titleLastIndex = html.index(titleLast, titleFirstIndex)
    print html[titleFirstIndex : titleLastIndex]
    findex = titleLastIndex
    
```

API 를 이용한 데이터 획득 (cf. Web crawling)



Big Data 의 전처리

1) 텍스트 데이터의 전처리: 자연어 처리



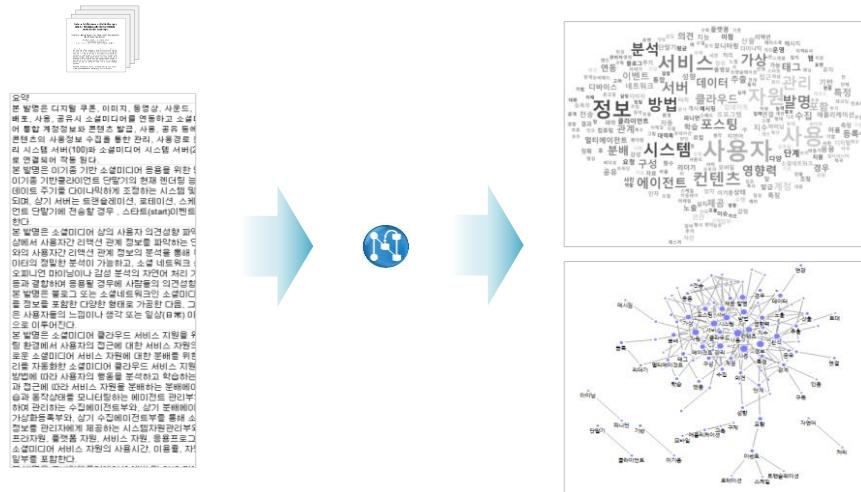
1) 텍스트 데이터의 전처리: 자연어 처리

Word Table			Word Occurrence Table		
	Part of Speech(POS)	Word length		등장한 문장	단어 등장 순서
국민	"Common Noun"	2	해결	41	315
국방력	"Common Noun"	3	추구	6	62
국정	"Common Noun"	2	정부	69	467
군림	"Common Noun"	2	평양	45	329
권력	"Common Noun"	2	기회	16	146
권위	"Common Noun"	2	경쟁	12	115
기관	"Common Noun"	2	어깨	31	257
기록	"Common Noun"	2	계층	66	448
기회	"Common Noun"	2	이웃	90	589
긴장	"Common Noun"	2	불행	74	492
길	"Common Noun"	1	일	85	566
끌	"Common Noun"	1	기관	38	296

추출된 단어, 품사, 글자수

단어가 문장/문단/문서에서 등장한 위치

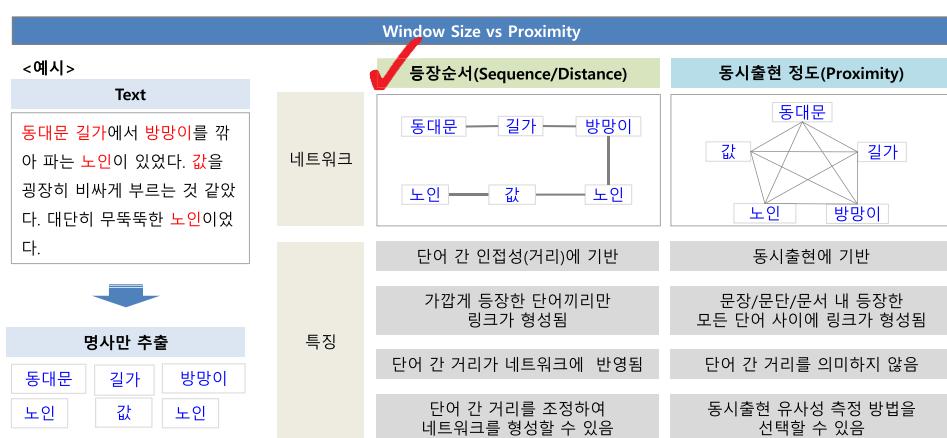
1) 텍스트 데이터의 전처리: Semantic network modeling



1) 텍스트 데이터의 전처리: Semantic network modeling

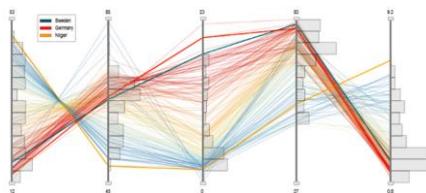
키워드 네트워크 구성

- 텍스트에서 추출한 단어들로 네트워크를 구성하는 방법은 두 가지로 구분할 수 있음
- 등장순서에 기반한 방법과 동시출현에 기반한 방법이 있음



2) Transaction data 의 전처리: network modeling

순번	보험종목	회사	사고발생일자	지급일자	병원이름	사유	피보험자주민번호	피보험자	계약자주민번호	계약자	수의자	설계사
1	생명보험	삼성생명	1998-12-30	1998-12-31	서울대병원	OO	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아나	나아나	최고임
2	생명보험	현대해상	1998-12-31	1999-01-01	연세대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아나	나아나	최고임
3	자동차보험	MG	1999-01-01	1999-01-02	순천향대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아나	나아나	최고임
4	상해보험	메리츠	1999-01-02	1999-01-03	서울대병원	DD	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아나	나아나	최고임
5	생명보험	삼성생명	2000-10-20	2000-10-20	연세대병원	AA	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아나	최고임
6	생명보험	현대해상	2000-10-20	2000-10-20	가천의대병원	BB	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
7	자동차보험	MG	2000-10-20	2000-10-20	서울대병원	CC	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
8	생명보험	삼성생명	2003-04-05	2003-04-05	연세대병원	EE	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
9	생명보험	현대해상	2003-04-05	2003-04-05	인하대병원	FF	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
10	자동차보험	MG	2003-04-05	2003-04-05	차병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
11	상해보험	메리츠	2003-04-05	2003-04-05	연세대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	나아나	최고임
12	생명보험	삼성생명	2010-02-03	2010-02-03	순천향대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아나	최고임
13	생명보험	삼성생명	2019-02-19	2019-02-19	서울대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
14	상해보험	메리츠	2019-02-19	2019-02-19	연세대병원	XX	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임



2) Transaction data 의 전처리: network modeling

순번	보험종목	회사	사고발생일자	지급일자	병원이름	사유	피보험자주민번호	피보험자	계약자주민번호	계약자	수의자	설계사
1	생명보험	삼성생명	1998-12-30	1998-12-31	서울대병원	OO	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아나	나아나	최고임
2	생명보험	현대해상	1998-12-31	1999-01-01	연세대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	나아나	나아나	최고임
3	자동차보험	MG	1999-01-01	1999-01-02	순천향대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아나	나아나	최고임
4	상해보험	메리츠	1999-01-02	1999-01-03	서울대병원	DD	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	나아나	나아나	최고임
5	생명보험	삼성생명	2000-10-20	2000-10-20	연세대병원	AA	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아나	최고임
6	생명보험	현대해상	2000-10-20	2000-10-20	가천의대병원	BB	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
7	자동차보험	MG	2000-10-20	2000-10-20	서울대병원	CC	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
8	생명보험	삼성생명	2003-04-05	2003-04-05	연세대병원	EE	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
9	생명보험	현대해상	2003-04-05	2003-04-05	인하대병원	FF	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
10	자동차보험	MG	2003-04-05	2003-04-05	차병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임
11	상해보험	메리츠	2003-04-05	2003-04-05	연세대병원	OO	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	나아나	최고임
12	생명보험	삼성생명	2010-02-03	2010-02-03	순천향대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	나아나	최고임
13	생명보험	삼성생명	2019-02-19	2019-02-19	서울대병원	XX	1985xxxx-2xxxxxx	홍길동	1946xxxx-1xxxxxx	김순이	김순이	최고임
14	상해보험	메리츠	2019-02-19	2019-02-19	연세대병원	XX	1946xxxx-1xxxxxx	김순이	1985xxxx-2xxxxxx	홍길동	홍길동	최고임

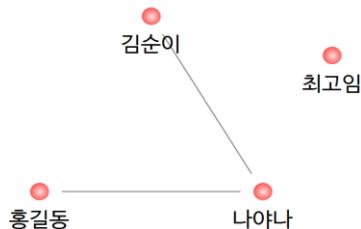
1 mode network



2) Transaction data 의 전처리: network modeling

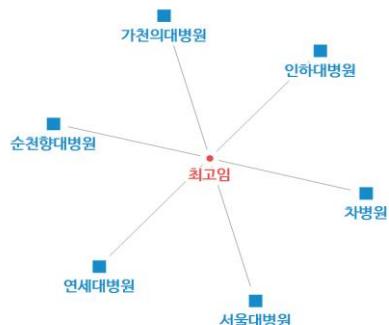
1 mode network

1-mode network Table(피보험자-계약자)

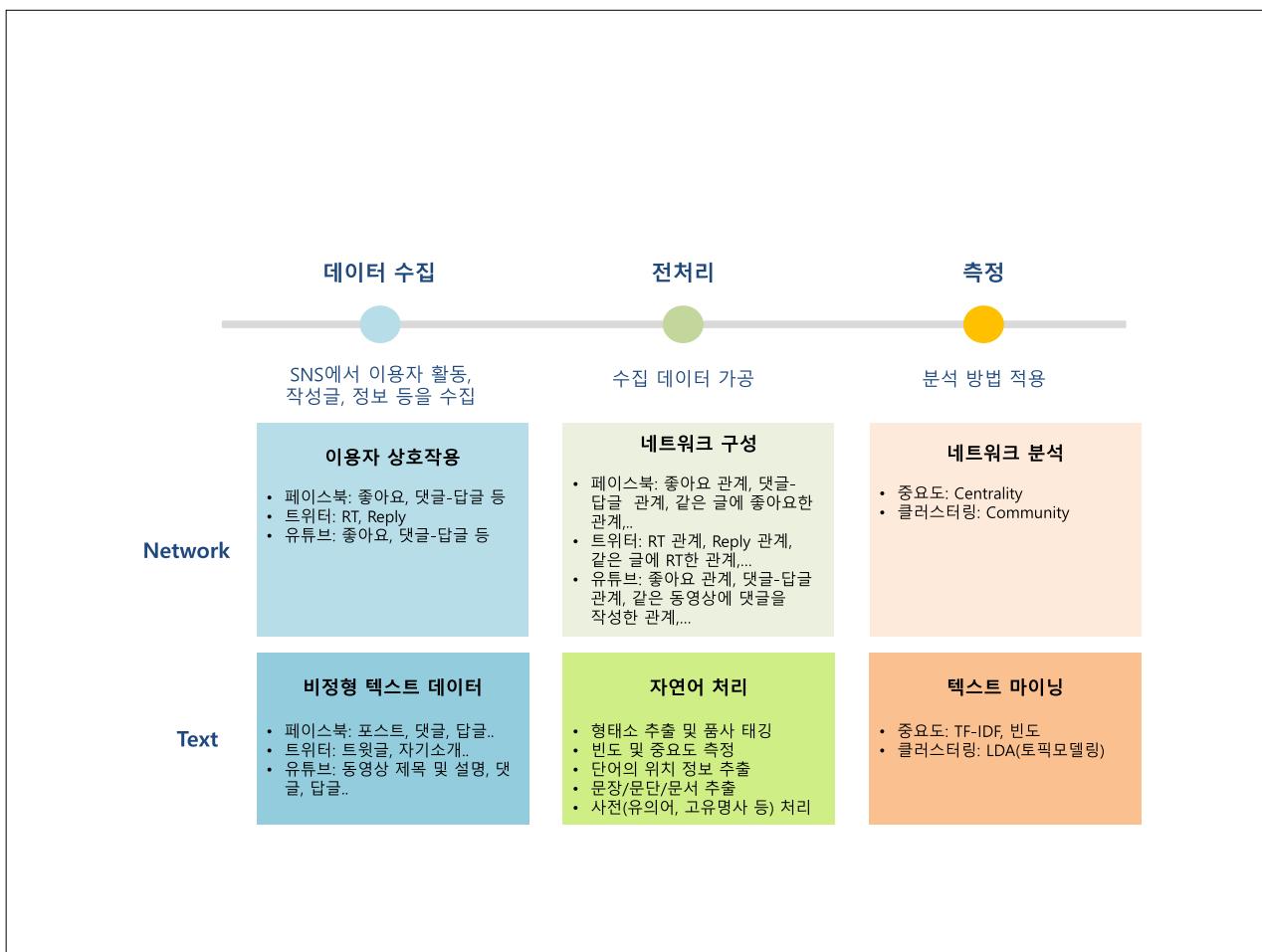


2 mode network

2-mode network Table(설계사-병원)



Big Data 의 분석



1) Text Mining

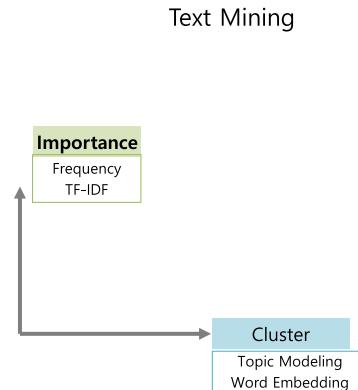
텍스트 데이터에 대한 체계적/정량적 분석의 필요성



	Structured Data	Unstructured Data
Quantitative Analysis	Statistics Social Network Analysis Machine Learning	???
Qualitative Analysis		Contents Analysis

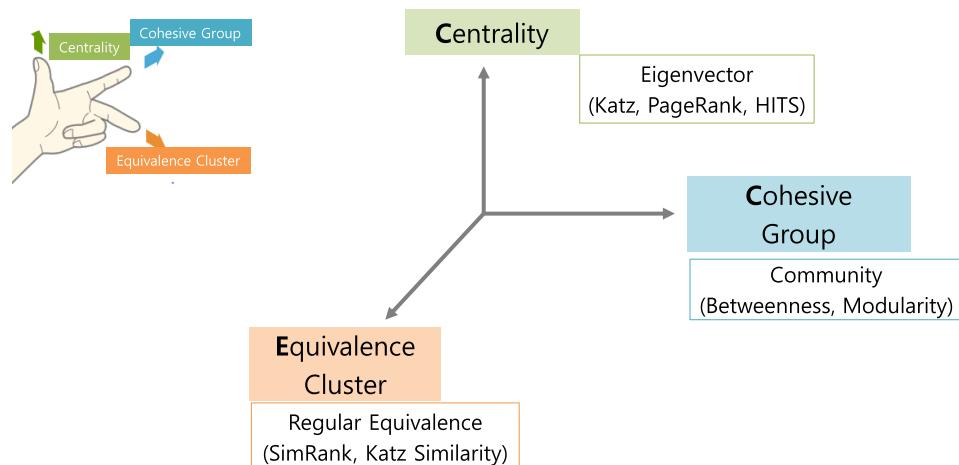
수작업과 분석자의 직관, 경험적 해석

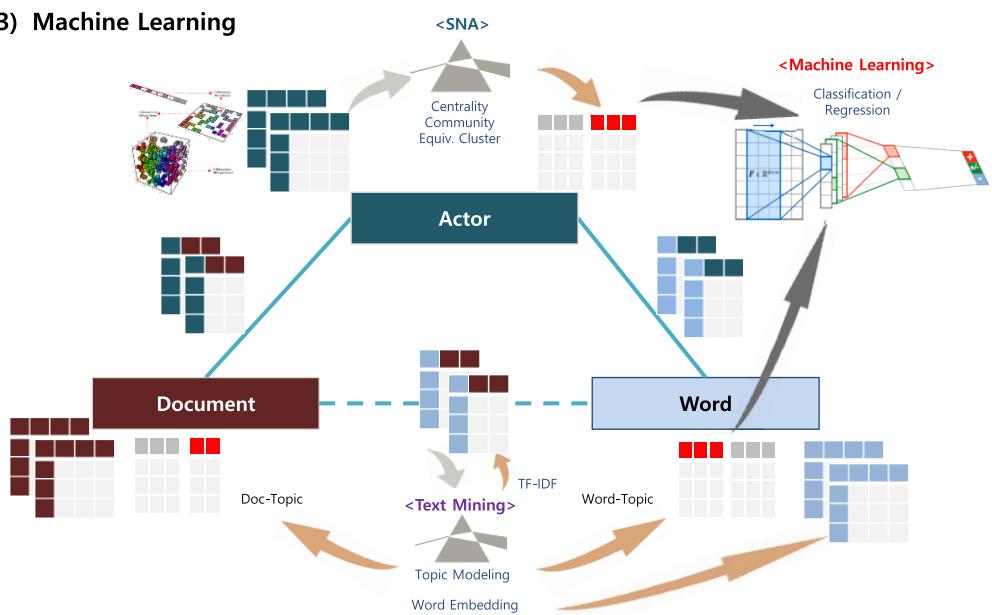
1) Text Mining



2) Social Network Analysis

Measure positions in social structure defined by intersecting relations



3) Machine Learning

사회과학 연구자를 위한 Big Data 교육과정

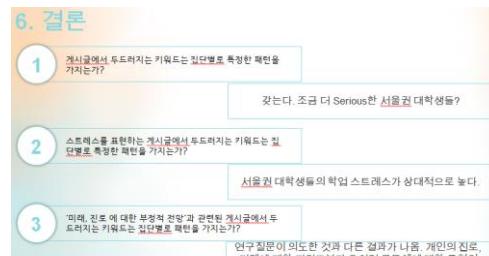
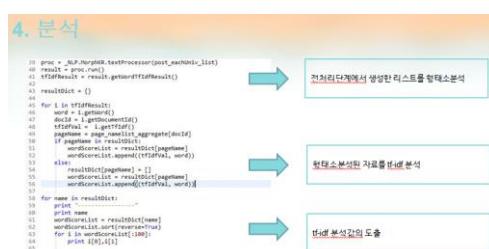
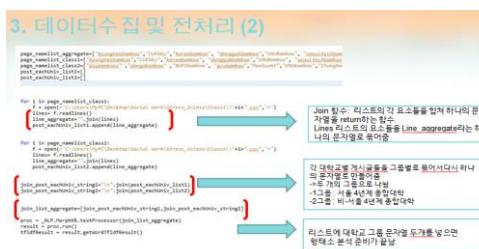
2019/1학 기/정규	학사	M130 4.0012 00	001	소셜 빅데이터 조사분석	전선	3-3-0	한국 어
-----------------	----	----------------------	-----	--------------	----	-------	---------

I. Python 문법

목차	세부 목차
1. 프로그래밍 입문(개요)	프로그래밍 언어 Why Python? 파이썬 설치 및 실행 파이썬 개발환경(IDLE)
2. 데이터	변수 Python의 자료형 (Data Type) Python의 자료구조 (Data Structure)
3. 흐름제어	조건문 반복문 예외처리
4. 함수	개요 Input Output Return 함수 안에서의 변수
5. 파일 입출력	파일 입출력의 필요성 파일 읽기 파일 쓰기
6. 클래스	클래스의 필요성 클래스의 정의 클래스의 사용
7. 모듈	모듈의 생성과 사용 내장 모듈 외부 모듈의 사용

II. 빅데이터 수집 분석 연습

목차	세부 목차
1. 빅데이터 분석 개요	빅데이터 입문 빅데이터 분석 절차
2. 수집	개발환경 크롤링, API, 인증 등
3) Web 크롤링	개요 게시판 제목 수집 네이버 실시간 검색어 수집 Tweepy 개요 트위터 가입 및 앱 생성
3) Twitter	Tweepy 설치 인증 데이터 수집 API 형태소 분석 데이터 구조화 NetMiner 입력
3. 전처리	통계 분석: 빈도 분석 핵심단어분석(TF-IDF) 소셜네트워크영향력 분석 단어-문서 연관성 분석 - word2vec
4. 분석	



개념어 검색 및 수집

수집, 전처리

```

searchList = ["#END_OF", "#DD", "#END_DD", "#END_EA", "#DE"]
tweetList = []

for keyword in searchList:
    userSet = set()
    userSet.add(keyword)
    print keyword
    lineString = ""
    line = 0

    tweepyCursor = tweepy.Cursor(api.search, q=keyword, count=100, result_type="recent", include_entities=True).items()

    while True:
        try:
            tweet = tweepyCursor.next()

            if line < 1000:
                lineString += tweet.text + "\n"
                line += 1
            else:
                tweetList.append(lineString)

                lineString = ""
                line = 0
        except:
            pass
    
```

Katz Status Centrality 분석

Retweet Network 수집, 전처리

```

if len(tweetList) >= 1000:
    userSet = set()
    userSet.add(user_id)

    try:
        seedUser = user.retweeted_status.user
        seedUser.id = seedUser.id
        if (seedUser.id, seedUser.id) in userNetwork:
            userNetwork[(seedUser.id, seedUser.id)] += 1
        else:
            userNetwork[(seedUser.id, seedUser.id)] = 1
    except:
        pass

    except:
        pass

    except:
        pass

    if len(userNetwork) >= 3000:
        break
    
```

NetMiner에 Network 자료 옮기기

```

    .to_csv('createNodeAll11c["node","user"]',index=False)
    .to_csv('createNodeAll11c["rela","rela"]',index=False)
    .to_csv('createNodeAll11c["rela","rela"]',index=False)
    .to_csv('createNodeAll11c["rela","rela"]',index=False)
    .to_csv('createNodeAll11c["rela","rela"]',index=False)
    
```

사회과학 연구자를 위한 Big Data 소프트웨어

개념어 리스트 간 Tf-Idf 분석

Tf-Idf 분석

```

proc = NLTKProcessor.tfidfProcessor(tweetList)
result = proc.run()
tfIdfResult = result.getTfIdfResult()

resultList = []
for i in tfIdfResult:
    docId = i.getDocID()
    word = i.getWord()
    searchKeyword = searchList[docID]
    if searchKeyword in resultList:
        wordCount = tfIdfResult[searchKeyword][word]
        wordScoreList.append((tfIdfVal, word))
    else:
        resultList.append(searchKeyword)
        wordCount = tfIdfResult[searchKeyword][word]
        wordScoreList.append((tfIdfVal, word))

for keyword in resultList:
    print keyword
    wordScoreList = result.getTfIdfResult[keyword]
    for i in wordScoreList[100]:
        print i[0], i[1]
    
```

상위 영향력자 트윗을 리스트 간 Tf-Idf 분석

수집 및 전처리

```

userSet = set()
tempList = []
for user in userList:
    userSet.add(user.id)
    print user.id
    tempList.append(user.id)
    for tweet in user.timeline:
        userTimeline.append(tweet.id)
        if len(userTimeline) >= 1000:
            break
    if len(userTimeline) >= 1000:
        break
    
```

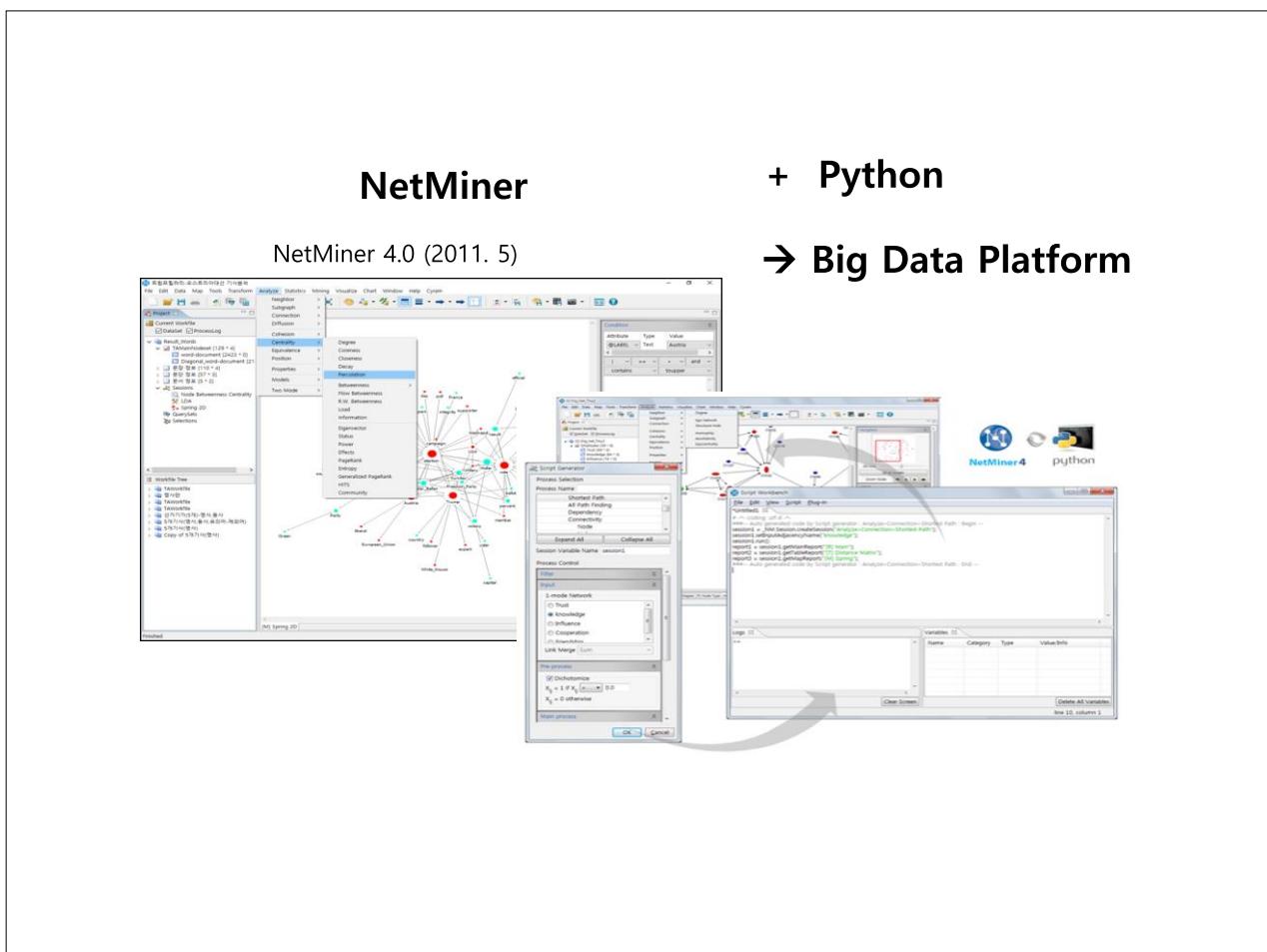
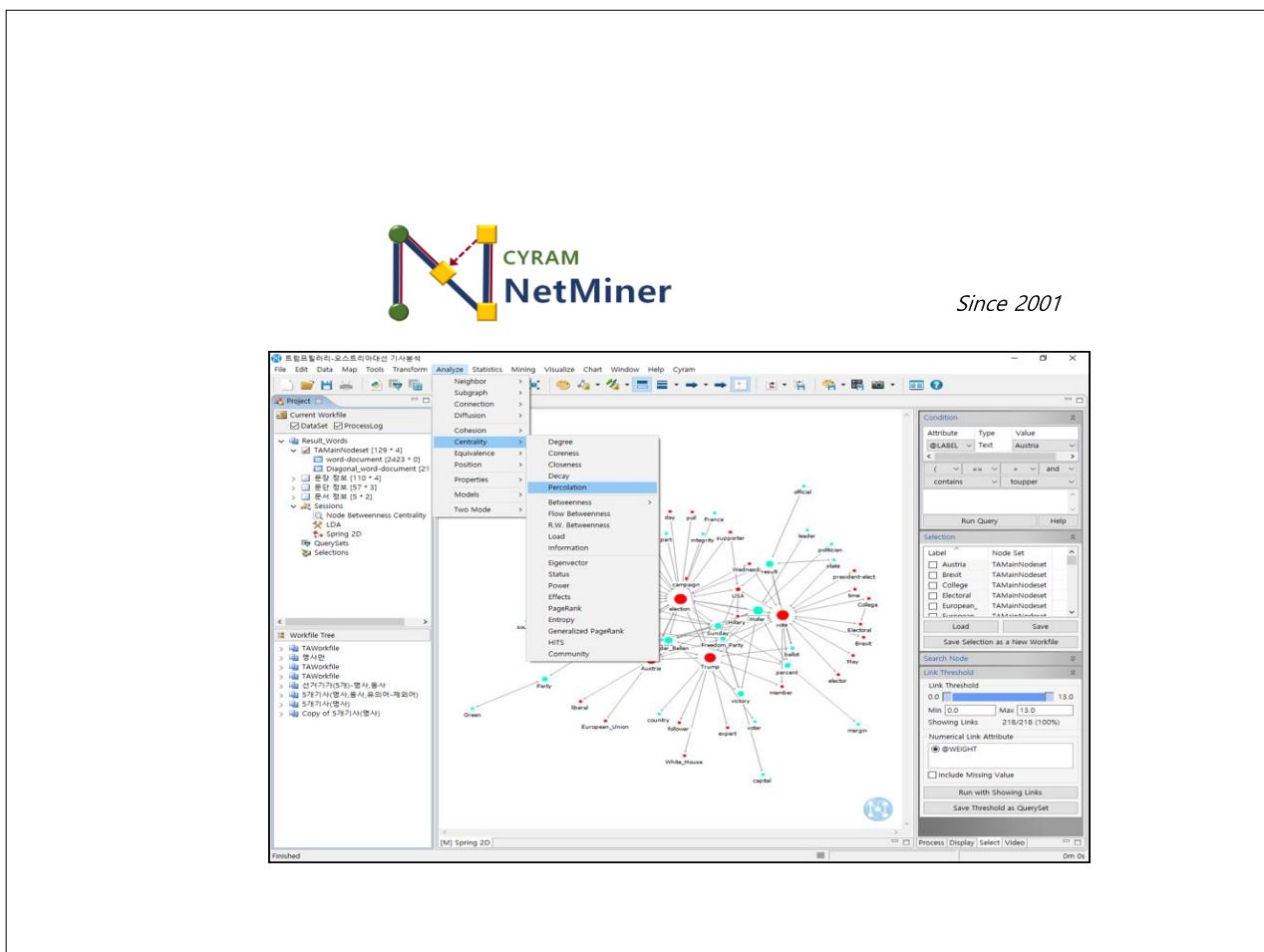
Tf-Idf 분석

```

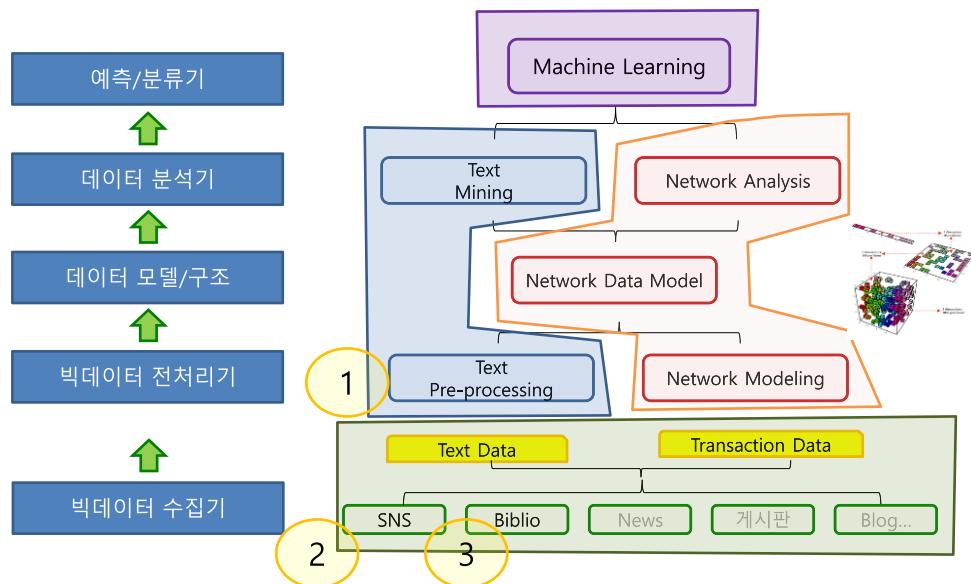
proc = NLTKProcessor.tfidfProcessor(userTimeline)
result = proc.run()
tfIdfResult = result.getTfIdfResult()

tfIdfList = []
for i in tfIdfResult:
    docId = i.getDocID()
    word = i.getWord()
    user = i.getUser()
    if user in tfIdfList:
        wordCount = tfIdfResult[user][word]
        tfIdfList.append((tfIdfVal, word))
    else:
        tfIdfList.append(user)
        wordCount = tfIdfResult[user][word]
        tfIdfList.append((tfIdfVal, word))

for user in tfIdfList:
    print user
    wordCount = tfIdfResult[user]
    for i in wordCount[100]:
        print i[0], i[1]
    
```

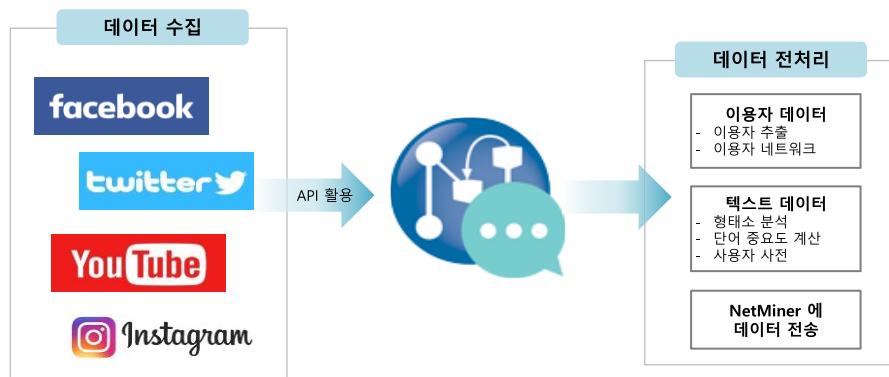


NetMiner 빅데이터 분석 Architecture



2. SNS Data Collector

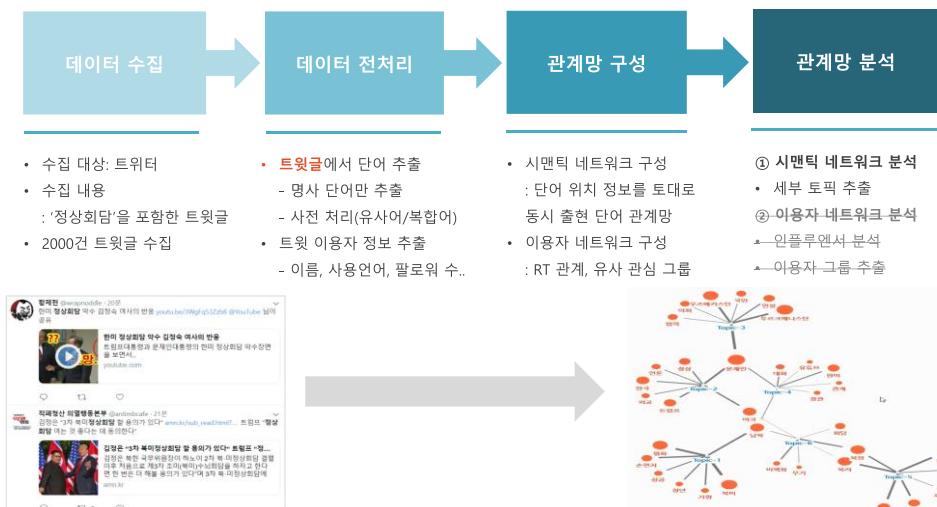
주요 소셜미디어의 공개 API(Application Programming Interface)를 이용하여 데이터를 수집하고,
수집한 데이터에 대한 전처리, 네트워크 모델링을 자동으로 수행하는
NetMiner의 확장 프로그램



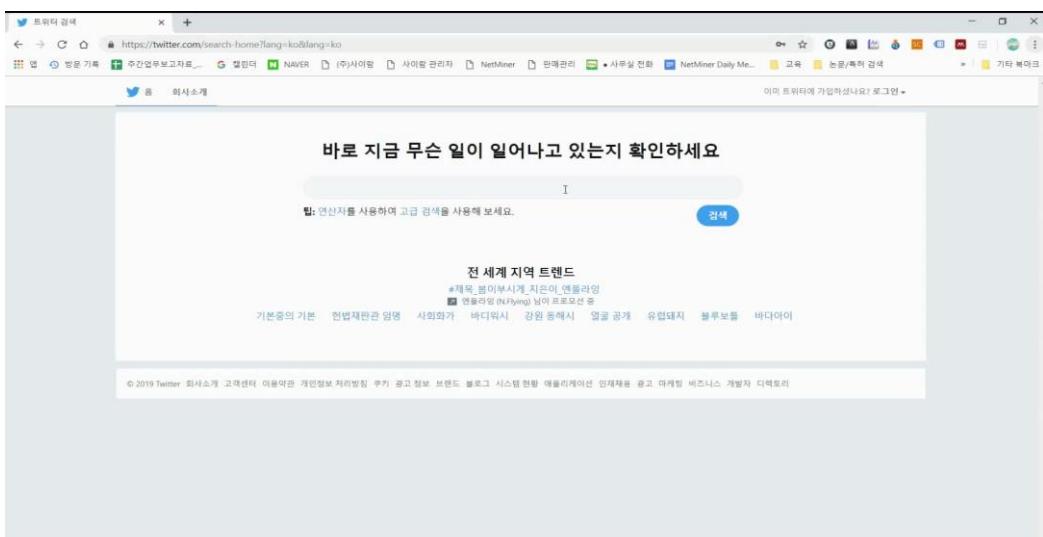
* Instagram은 5월 출시 예정

2. SNS Data Collector

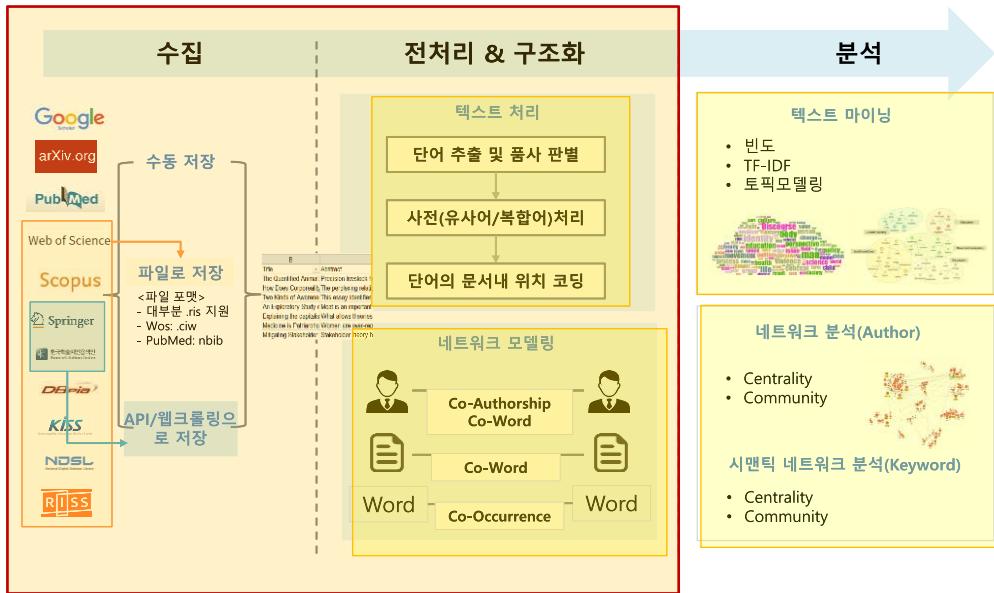
- 수집데이터: '정상회담' 관련 트윗
- 수집기간: 19년 4월 3째주(2019.4.13~4.19)



NetMiner SNS Data Collector

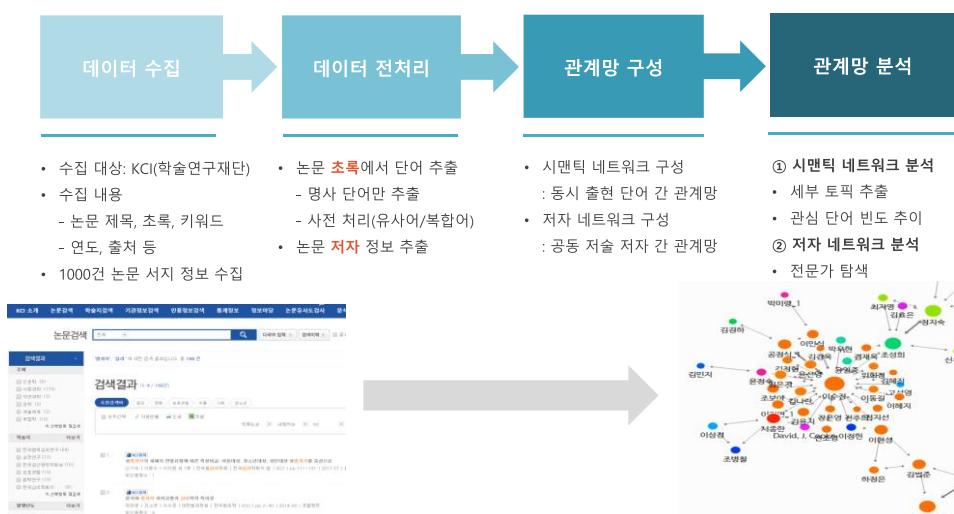


3. Biblio Data Collector

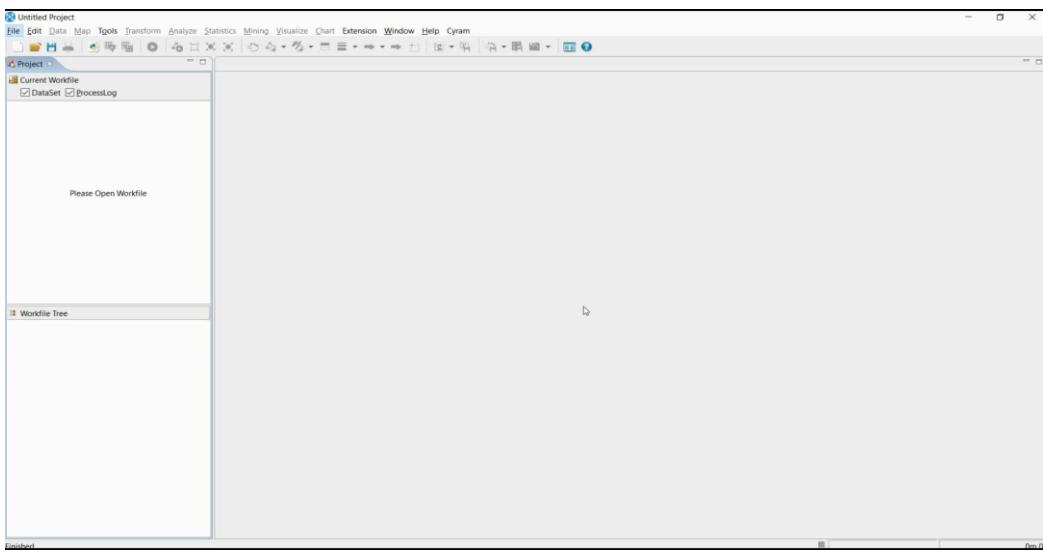


3. Biblio Data Collector

- 수집데이터: 범죄자 심리를 연구한 국내 논문 서지 정보(문헌 정보)
- 수집기간: 2000.1~2019.3

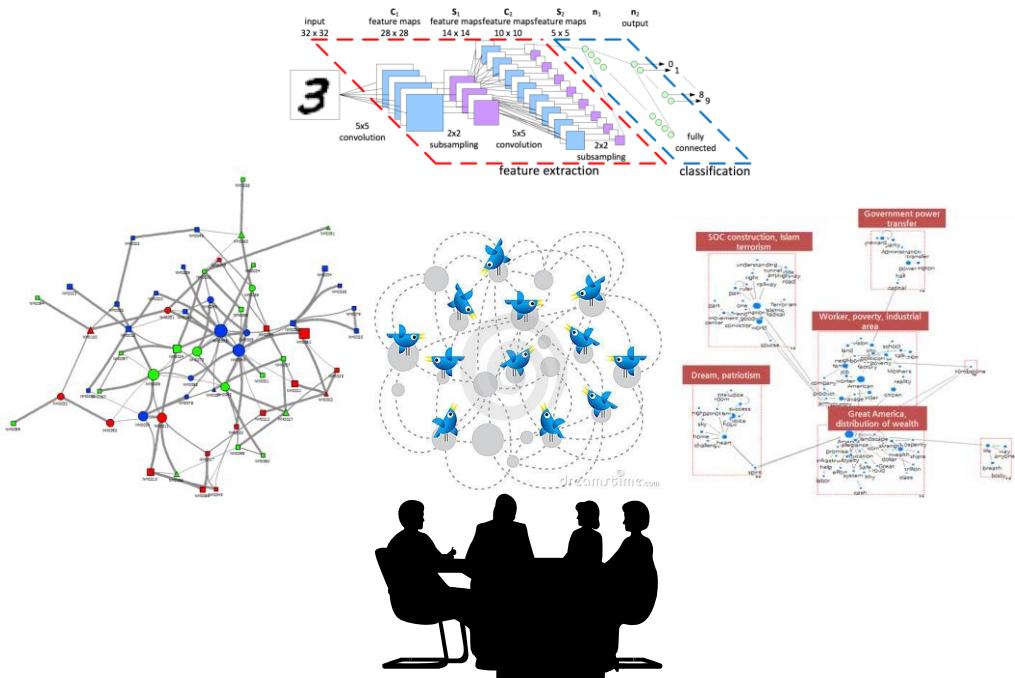
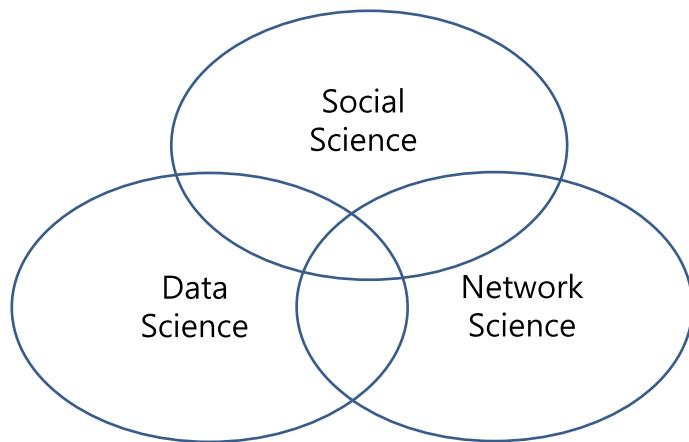


NetMiner Biblio Data Collector



맺음말

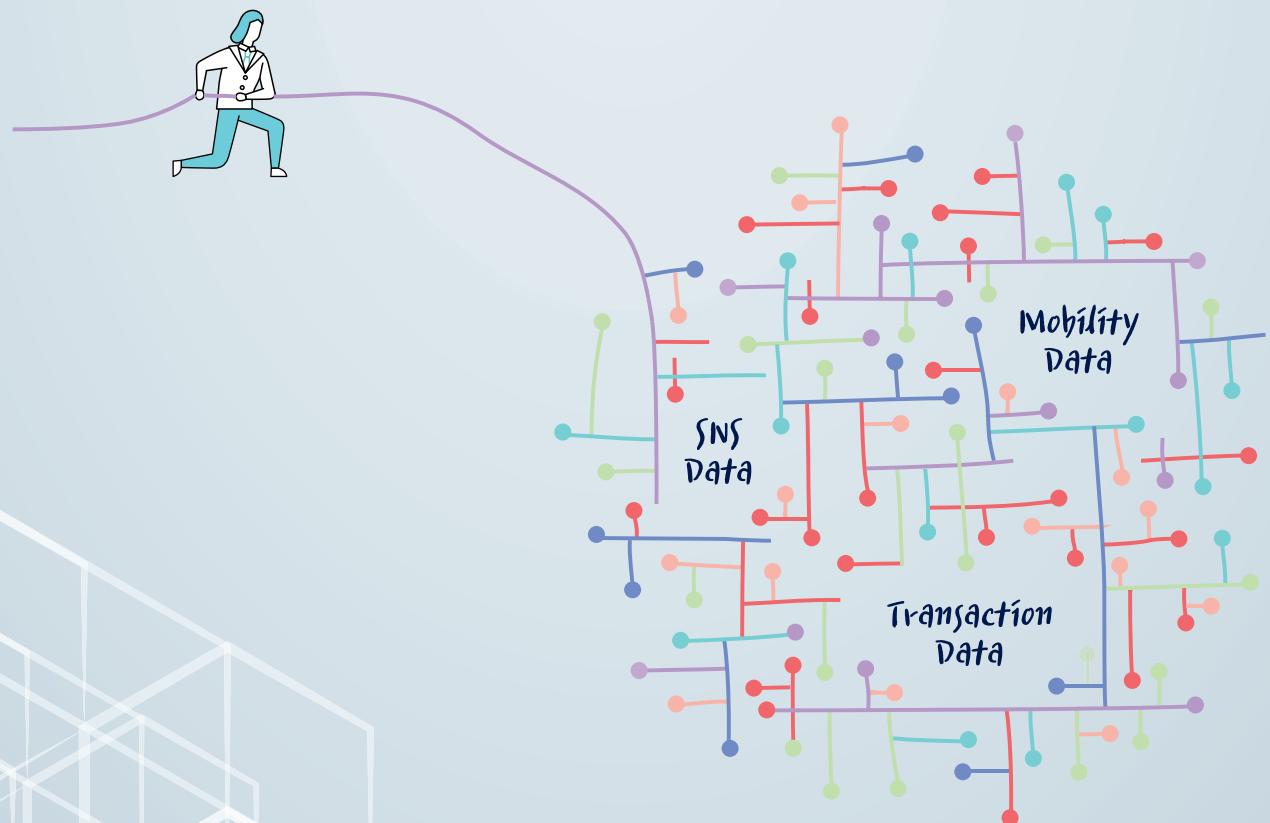
Future direction of social sciences ?



▶ 2부: 빅데이터 활용 사례

모빌리티 데이터로 바라보는 사회

김정민 연구원 (카카오 모빌리티)



모빌리티 데이터로 바라보는 사회

김정민 / 데이터랩 / 카카오모빌리티

dominic.jmkim@kakaomobility.com

KOSSDA 데이터페어

2019. 06. 27.

1

Contents

- 모빌리티 데이터
- 데이터 분석 방법론
- 모빌리티 데이터로 바라보는 사회 (연구사례)

2

모빌리티 데이터

3

모빌리티 데이터



승객



기사



호출



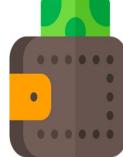
운행



위치



지도



결제



프로모션

수요공급 불일치 문제를 완화하여 승객과 기사 모두에서 '빠르고 편리한 운행'을 제공하는 것이 목표

4

모빌리티 데이터

장소(Place) = 위치(Location) + 맥락정보(Context)

위치정보 변환: 단말 신호 (GPS, Wifi, LTE 등)

맥락정보: 카카오맵에서 장소와 카테고리를 지속적으로 관리

** 주소, 상호, 카테고리 등

데이터	장점	단점
단순 휴대폰 위치 데이터	지속적으로 정확한 위치 파악 가능	사용자 방문지, 혹은 행위 유추가 어려움
금융 결제 데이터	정확한 행위 파악 가능	이동 경로, 결제가 없는 행위 유추 불가
대중교통 승하차 데이터	이동 경로 및 시간 정확히 파악 가능	행위 유추 및 최종 목적지 확인 불가
검색 데이터 [다음, 네이버 등]	관심사 파악 가능, 행위 예측 가능	실제 행위 여부 파악 불가
모빌리티 이동 데이터	이동 경로 및 시간, 비결제를 포함한 행위 파악 가능	사용하지 않는 시간 대가 존재

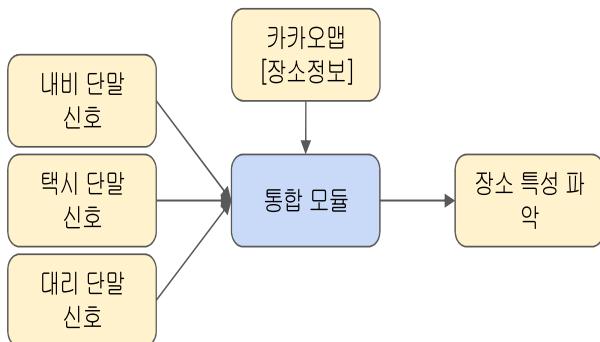
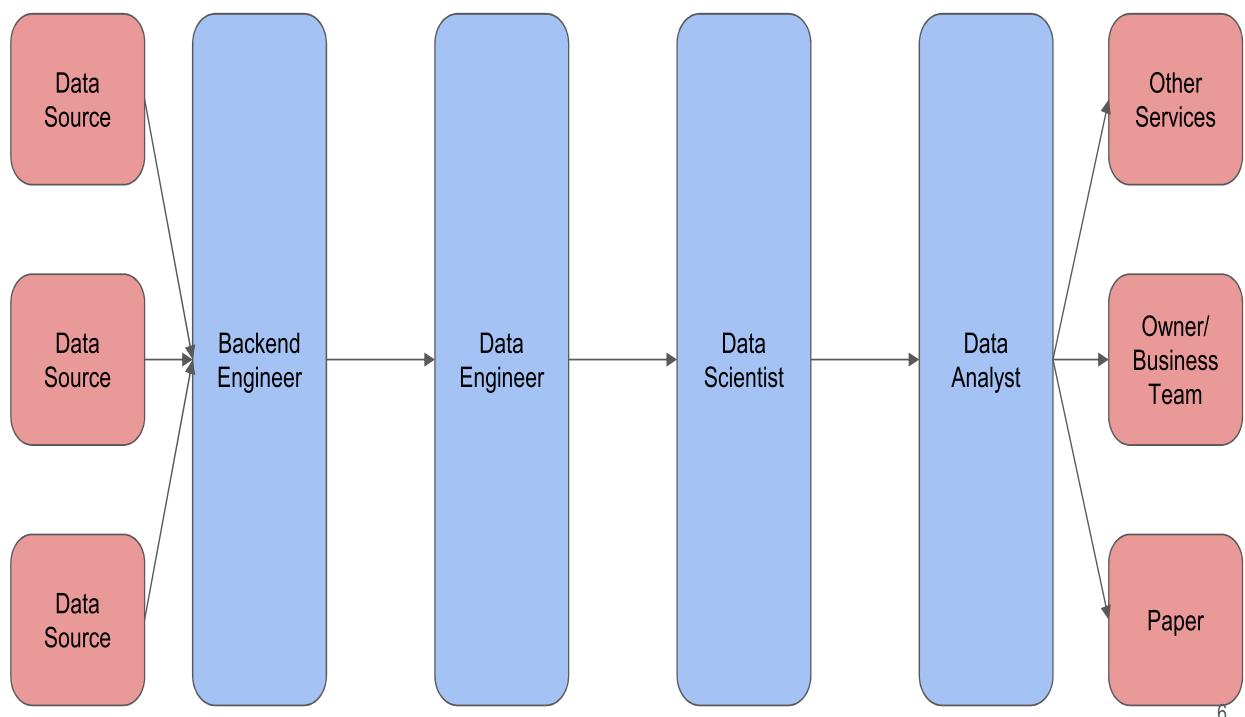


표: 사용자 위치 및 행동 파악을 위한 데이터 비교

5

데이터 분석 흐름



6

ETL (Extract, Transform, Load)

데이터는 그 의미를 추출(Extract)하고 분석 가능한 형태로 변환(Transform)하여 적재(Load)해야 한다.

추출: 텍스트(자연어처리), 이미지(물체 인식), 사운드 (주파수 분석) 등..

변환: 분석 하고자 하는 내용을 Byte/Number/String, 정형/비정형 등 적절한 형태로 변환

** 추출, 변환 및 저장의 효율을 높이기 위해서 다양한 분산처리 시스템 (Distributed System)이 사용된다. [Hadoop 등]

적재 방식:

Batch → Micro Batch → Streaming

Batch: 주기적으로 ETL 작업을 실행 (매달, 매주, 매일, 매시 ..)

Micro Batch: Batch의 주기를 매우 짧게하여 Streaming과 유사하게 처리

Streaming: 데이터를 하나의 흐름처럼 실시간으로 처리



Streaming Framework

저장소: DB(SQL/NoSQL), File System(Local/HDFS)



SQL DB



NoSQL DB



분산처리 프레임워크
HDFS (Hadoop File System)
MapReduce: 병렬처리 연산 프레임워크

7

대한민국 모빌리티 연구의 현실?



학교, 연구소

연구할 사람은 많은데
좋은 데이터가 없다



회사

데이터는 많은데
연구할 사람이 없다

VS

8

데이터 분석 방법론

9

시계열 분석

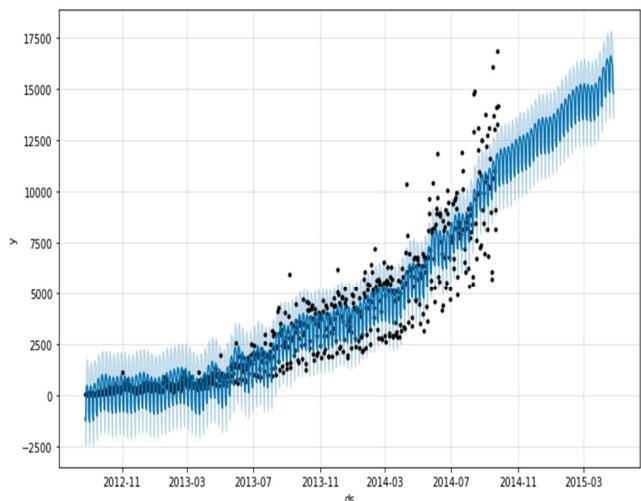
시계열 분석은 시간대에 따른 값을 통해
패턴이나 특이점을 찾아내는 분석

경향(trend): 장기적 변화

계절성(seasonality): 시간대별 유사 패턴

주기(cycle): 일정한 진폭마다 유사 변동

불규칙성(irregular): 오차, 잡음, 이벤트



PROPHET

10

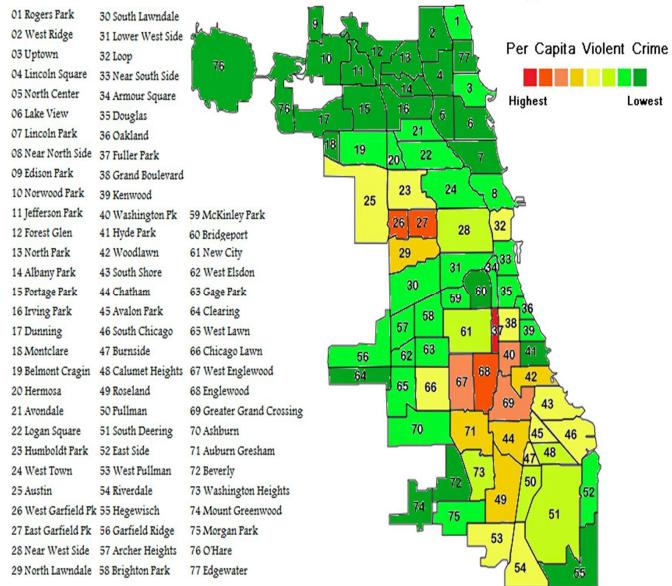
공간 통계 분석

공간 통계 분석은 위상, 기하, 지리적 특성을 사용하여 개체를 설명하는 방법이다.

지리적으로 가까울수록 유사한 특성을 나타낸다.
(First Law of Geography, Tobler, 1970)

→ 공간 자기 상관 (Spatial Autocorrelation)

일반적으로 통계에서는 각 사건의 독립성을 전제한다.
그러나 공간 통계에서는 지리적 유사성, 즉 공간 자기상관을 고려하여 분석한다.



11

네트워크 분석

Point/Node: 하나의 개체

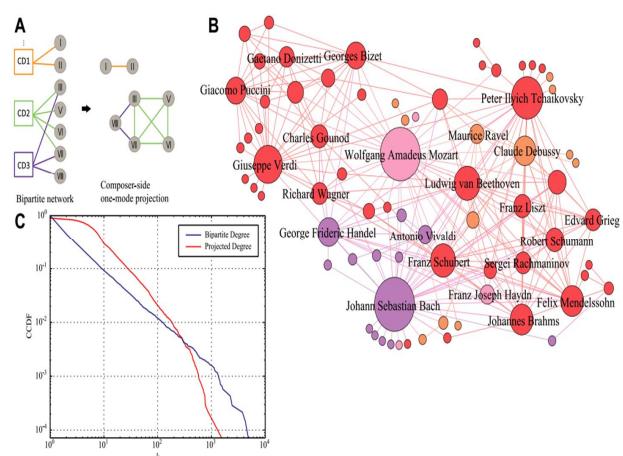
Link/Edge: 개체 사이의 관계

분석 방법:

중심도(centrality), 밀도(density),
반경(diameter) ..

주 분석 분야:

사회관계망, 인프라망, 문서분석 ..

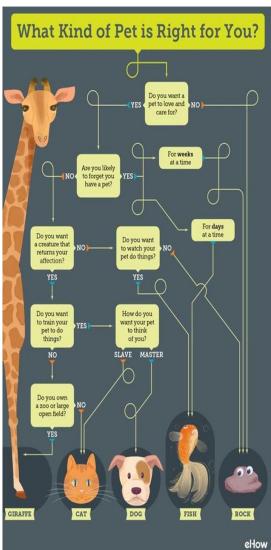


Park, D., Bae, A., Schich, M. et al. EPJ Data Sci. (2015) 4: 2. <https://doi.org/10.1140/epjds/s13688-015-0039-z>
[고전 음악가 네트워크 분석 - 발표자 주]

12

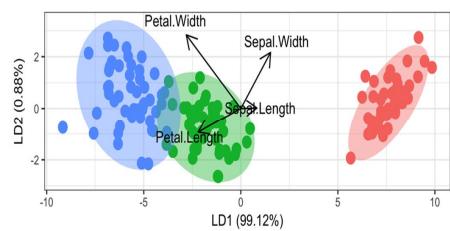
분류 분석 (classification)

특정 데이터들을 몇 개의 정해진 기준으로 분류하는 분석



<https://www.pinterest.com/pin/76420524906163543/>

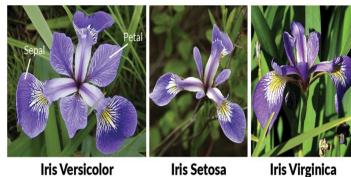
미리 정의된 그룹으로 데이터를 분류
: Decision Tree, Neural Networks,
SVM (Support Vector Machine), Ensemble
Method, kNN (k Nearest Neighborhood),
Logistic Regression ..



<https://lchblogs.netlify.com/post/2017-12-22-reddconellipsida/>

데이터 내부 특성을 사용하여 자동 분류

: Clustering



<http://www.lac.inpe.br/~rafael.santos/Docs/R/CAP394/WholeStory-Iris.html>

13

자연어 분석

인간의 언어를 기계를 이용해서 분석

기계학습의 하위 분야로 인식

분석 순서(한국어): 형태소 분석 → 품사 부착
→ 구절단위분석 → 구문분석

정보 검색, 정보 추출, 음성 인식, 단어분류,
문장/문서 분류, 감정분석, 기계번역



Figure 4: Anger.

Figure 5: Fear.

Figure 6: Disgust.



Figure 7: Happiness.

Figure 8: Sadness.

Figure 9: Surprise.

소셜 미디어 기반(geo tagged tweets)의 감정 지도 (워싱턴) [발표자 -주]

Yerach Doytsher, Ben Galon, and Yaron Kanza. 2017. Emotion Maps based on Geotagged Posts in the Social Media. In Proceedings of the 1st ACM SIGSPATIAL Workshop on Geospatial Humanities (GeoHumanities'17). ACM, New York, NY, USA, 39-46. DOI: <https://doi.org/10.1145/3149858.3149862>



14

이미지 분석

이미지 분석은 인간이 컴퓨터보다 더 잘하는 분야이다.

이미지의 맥락을 이해하는 것은 오랜 기간 어려운 문제였다.

Google Cloud Vision API:

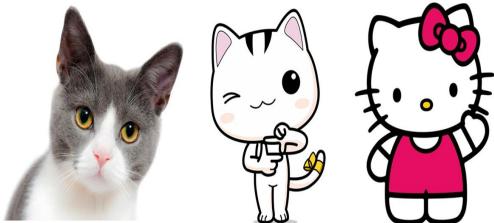
Label Detection(사물 인식)

Logo Detection(회사 로고 인식)

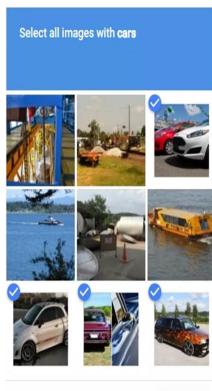
Landmark Detection (랜드마크 인식)

Face Detection (사람 얼굴 찾기 및 감정 분석)

Optical Character Recognition (문자인식) 등



위의 이미지들은 고양이인가요?



reCAPTCHA

매크로 방지 및 이미지 분석 태그 획득

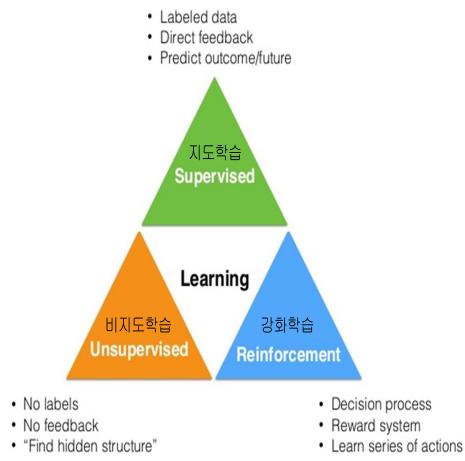
고양이 사진: <https://www.wadiz.kr/web/campaign/detail/20802>

고양고양이: <https://www.facebook.com/goyangcity>

헬로키티: https://www.altpress.com/news/hello_kitty_creators_reveal_character_isnt_actually_a_cat/

15

머신 러닝



지도 학습(Supervised Learning): 데이터에 대한 명시적인 정답이 주어진 상태에서 컴퓨터를 학습 / classification, regression

비지도 학습(Unsupervised Learning): 데이터에 대한 명시적인 정답이 주어지지 상태에서 컴퓨터를 학습 / clustering

강화 학습(Reinforcement Learning): 에이전트가 주어진 환경(state)에 대해 어떤 행동(action)을 취하고 이로부터 어떤 보상(reward)을 얻으면서 학습, 에이전트는 보상(reward)을 최대화(maximize)하도록 학습 / alphago

많은 경우에 머신러닝은 기존에 있던 분석 방법들을 더 빠르고 정확하게 수행해 주는 역할을 한다.

16

모빌리티 데이터로 바라보는 사회

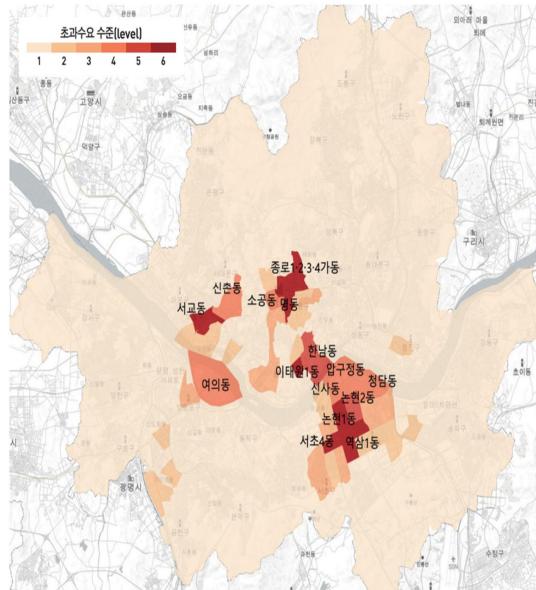
17

시공간 분석

18

시야 교통 불균형 현황(택시) 빅데이터 분석

<그림 4-3> 서울의 심야시간 택시 초과수요 현황



<표 4-1> 심야시간 택시의 초과수요가 많은 지역

초과수요 수준	초과수요가 많이 발생한 행정동
Level 6-2	역삼1동(강남구), 종로1-2-3-4가동(종로구)
Level 6-1	서교동(마포구), 논현1동(강남구), 명동(중구), 서초4동(서초구), 이태원1동(용산구)
Level 5	논현2동(강남구), 한남동(용산구), 여의동(영등포구), 압구정동(강남구), 청담동(강남구), 신사동(강남구), 신촌동(서대문구), 소공동(중구)

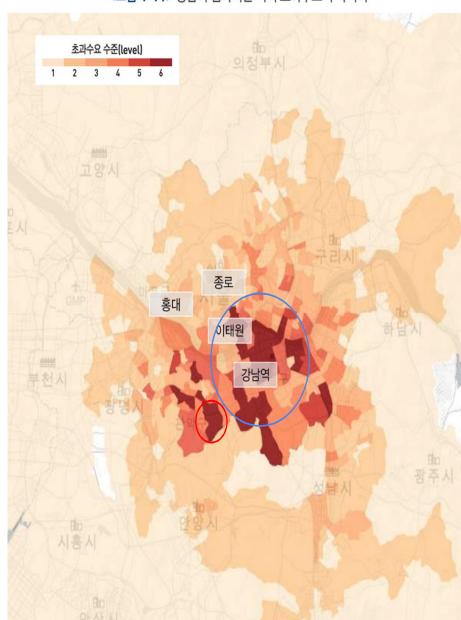
초과수요:

특정 시간대에 택시를 이용하고자 했으나 이용하지 못한 승객(혹은 콜) 수

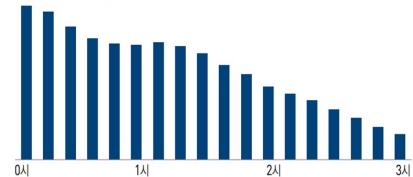
시민 이동성 증진을 위한 심야 교통 현황 분석 [카카오모빌리티, 서울디지털재단, 2018. 12.]

시야 교통 불균형 현황(택시) 빅데이터 분석

<그림 4-11> 강남역 심야시간 택시 초과수요의 목적지



<그림 4-10> 강남역 심야시간 시간대별 택시 초과수요 분포

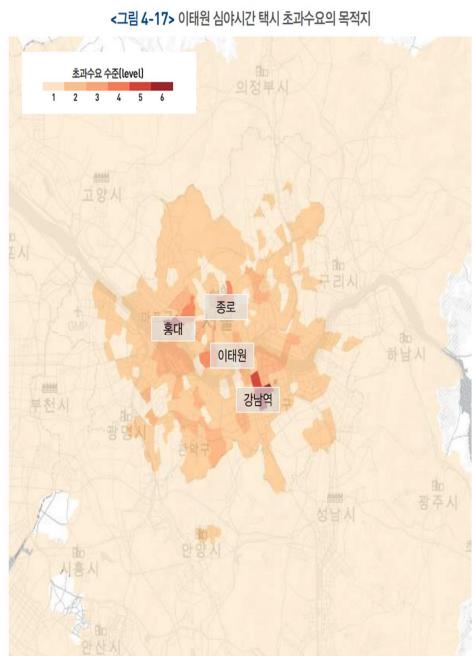


<표 4-3> 강남역 심야시간 택시 초과수요의 목적지

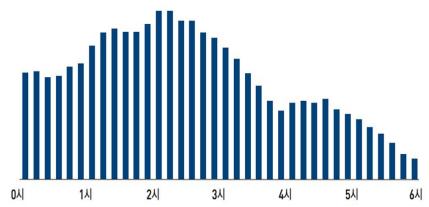
초과수요 수준	초과수요의 목적지(행정동)
Level 6-2	논현1동(강남구), 역삼1동(강남구), 청담동(강남구), 이태원1동(용산구), 잠원동(서초구), 서원동(관악구), 행운동(관악구), 논현2동(강남구), 압구정동(강남구)
Level 6-1	삼성2동(강남구), 신림동(관악구), 인화동(관악구), 역삼2동(강남구), 지양4동(광진구), 청량리동(관악구), 남현동(관악구), 문장2동(송파구), 서초2동(서초구), 상도1동(동작구), 양재1동(서초구), 잠실2동(송파구), 한남동(용산구), 잠실본동(송파구), 양재2동(서초구)
Level 5	대치4동(강남구), 회현동(광진구), 잠실3동(송파구), 세곡동(강남구), 대방동(동작구), 삼전동(송파구), 반포동(서초구), 석촌동(송파구), 서초4동(서초구)
Level 4	사당1동(동작구), 복지2동(송파구), 신사동(강남구), 삼성1동(강남구), 흑석동(동작구), 대치2동(강남구), 서강2동(동작구), 행당1동(성동구), 자양2동(강진구), 구로3동(구로구), 노량진1동(동작구), 양재1동(동대문구), 구의3동(용산구), 대학동(관악구), 신사동(관악구), 옥수동(성동구)

시민 이동성 증진을 위한 심야 교통 현황 분석 [카카오모빌리티, 서울디지털재단, 2018. 12.]

시야 교통 불균형 현황(택시) 빅데이터 분석



<그림 4-16> 이태원 심야시간 시간대별 택시 초과수요 분포



<표 4-6> 이태원 심야시간 택시 초과수요의 목적지

초과수요 수준	초과수요의 목적지(행정동)
Level 6-1	서교동[마포구], 역삼1동[강남구]
Level 5	논현1동[강남구]
Level 4	없음

시민 이동성 증진을 위한 심야 교통 현황 분석 [카카오모빌리티, 서울디지털재단, 2018. 12.]

21

Abnormal(Event) Detection

22

이동 패턴 분석 [핫스팟 분석]

열대야 기간 어디 많이 갔나

단위:급상승 횟수, 건, ()안은 %

실내

대형마트·몰	699(21.2)
멀티플렉스 극장	264(8.0)
아쿠아리움, 박물관 등	145(4.4)
찜질방	17(0.5)

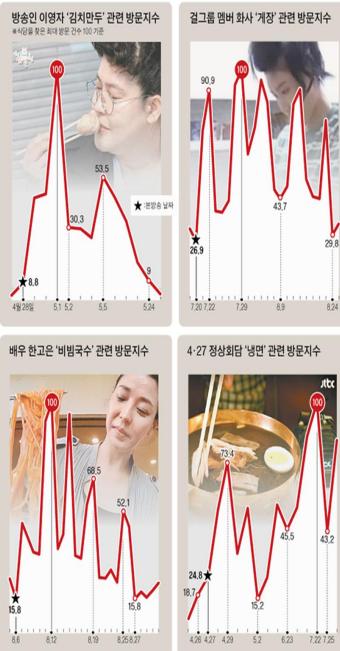
야외

해수욕장, 워터파크	479(14.5)
테마파크·공원	297(9.0)
계곡	135(4.1)
동굴·폐тан광	81(2.5)

기타

맛집	277(8.4)
호텔	141(4.3)
병원·장례·예식·종교	128(3.9)
전통시장	79(2.4)
주거·교통 등	554(16.8)

*열대야 : 밤 최저 기온이 섭씨 25도를 웃도는 기간
(올해 7월 22일~8월 16일).



시계열 분석을 통해 평소와 다른 이상패턴을 찾아낼 수 있고 이를 기반으로 사람들의 생활 패턴 변화를 살펴볼 수 있다.

단순히 많이 간 곳으로 분석하면 언제나 1위는 비슷하다.

평소보다 갑자기 더 많이 방문한 곳 혹은 카테고리를 분석해야 패턴의 변화를 살펴 볼 수 있다.

23

이동 패턴 분석 [핫카테고리 분석]

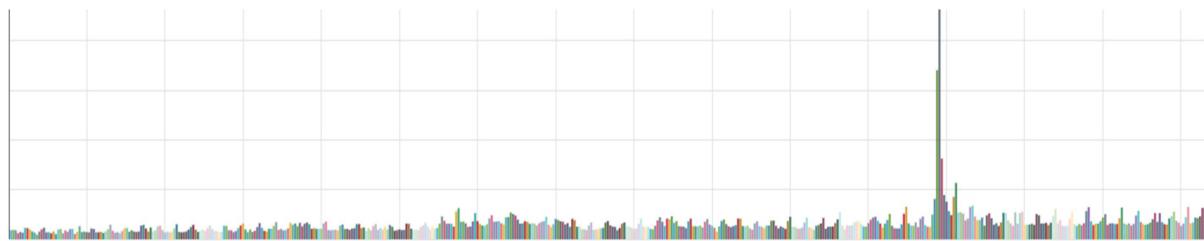
1월 26일부터 약 1주간 빨래방의 급증!?

얼어붙은 세탁기에 빨래방 '한파 특수'...동파 막으려면?

(http://news.jtbc.joins.com/article/article.aspx?news_id=NB11581868)

빨래방 카테고리 일간 길안내 수 (2017. 01 ~ 2018. 04)

2018. 01. 28.



24

이동 패턴 분석 [핫카테고리 분석]

2018년 설 기간 핫 카테고리 분석

4일	5일	6일	7일	8일	9일	10일
					2018 동계 올...	하나로클럽, 떡/한과, 국립묘지, 초콜릿
11일	12일	13일	14일	15일	16일	17일
하나로클럽, 납골당, 국립묘지, 장례, 묘지, 정관장, 초콜릿	하나로클럽, 전통식 품제조, 정관장	초콜릿, 정관장, 하 나로클럽, 도축업, 정육점, 떡/한과	정관장, 초콜릿, 떡/ 한과, 하나로클럽, 과일, 채소가게	설날 연휴 방앗간, 제시음식, 마을회관, 세차장	설날 요양원, 장례, 납골 당, 마을회관, 화장 터, 빌라, 아파트	설날 연휴 요양원, 장례, 납골 당, 묘지, 스케이트 장
18일	19일	20일	21일	22일	23일	24일
장례, 스케이트장, 납골당	스케이트장, 입시학 원					
25일	26일	27일	28일	3월 1일	2일	3일
2018 동계 올...				삼일절		

25

Clustering

26

Living Zone Analysis

- 행정동 간 통행량을 기반으로 네트워크를 생성하고, 클러스터링을 통해서 생활권을 구분
 - ** 이 작업에서 물리적 위치, 거리 정보는 별도로 사용하지 않음
 - ** Fast greedy modularity optimization algorithm 사용
- 행정동 간 통행량은 물리적 거리와 매우 높은 상관관계를 보이기 때문에 인접한 곳들이 하나의 클러스터, 즉 생활권으로 표현됨
- 주요 분석 내용
 - 각 생활권의 중심 지역(허브)
 - 각 생활권 내부에서의 주요 통행
 - 시간대에 따른 생활권의 변화
- 생활권 분석은 수요/공급 모델의 정교화 및 공공 정책 수립 근거자료 등으로 사용 가능

27

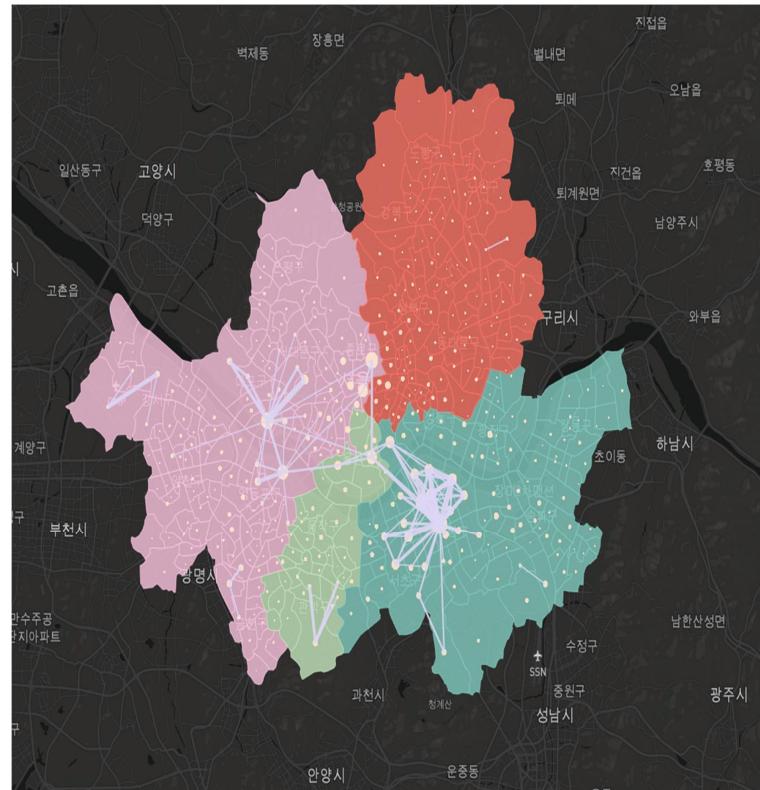
Living Zone Analysis

서울/모든 요일/모든시간
분석 기간: 20180901~20181030
분석데이터: 카카오모빌리티 택시 출도착

4개 생활권(임의 명명):
동남권, 동북권, 서부권, 관악권

주요 거점:
강남, 흥대, 광화문, 서울대

서울대를 중심으로 하는 관악권이 별도의 생활권으로 분리되는 이유는?



28

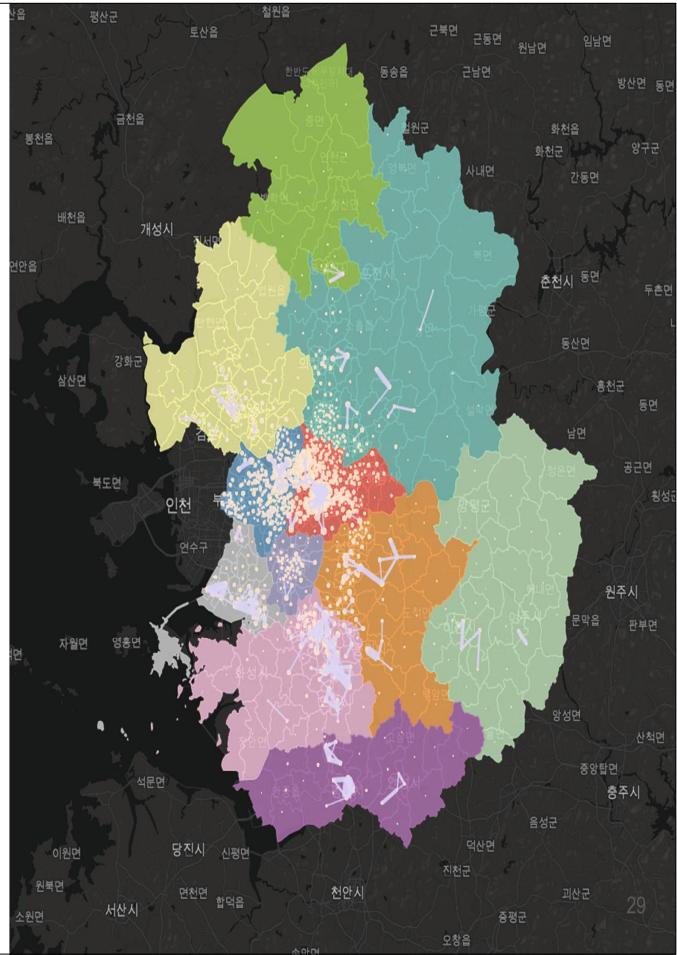
Living Zone Analysis

서울, 경기/모든 요일/모든 시간

분석 기간: 2018.09.01 ~ 2018.10.30

11개 생활권:

- 고양-김포-파주권: 일산 중심 통행 많음
- 서울서부권
- 서울동부권
- 서울동북-의정부-양주-포천권
- 성남-용인-광주권: 광주 - 성남 간 통행 많음
- 수원-화성-오산: 내부 통행 많음
- 안양권
- 안산-시흥권: 내부 통행 많음
- 여주-이천권
- 연천권
- 평택권: 평택, 송탄 근방 내부 통행 많음



Summary

- 모빌리티 데이터: 우리의 삶을 보여주는 중요한 데이터 중 하나
- 모빌리티 데이터 분석
 - 모빌리티 데이터 + 서비스 데이터
 - 다양한 방법론 적용
- 모빌리티 데이터 분석의 결과
 - 보다 빠르고 정확하고 편리한 이동 제공
 - 이동 최적화를 통한 사회적 비용 감소
 - 각 사용자에게 적합한 서비스 제공
 - 행위에 대한 다양한 가설을 검증

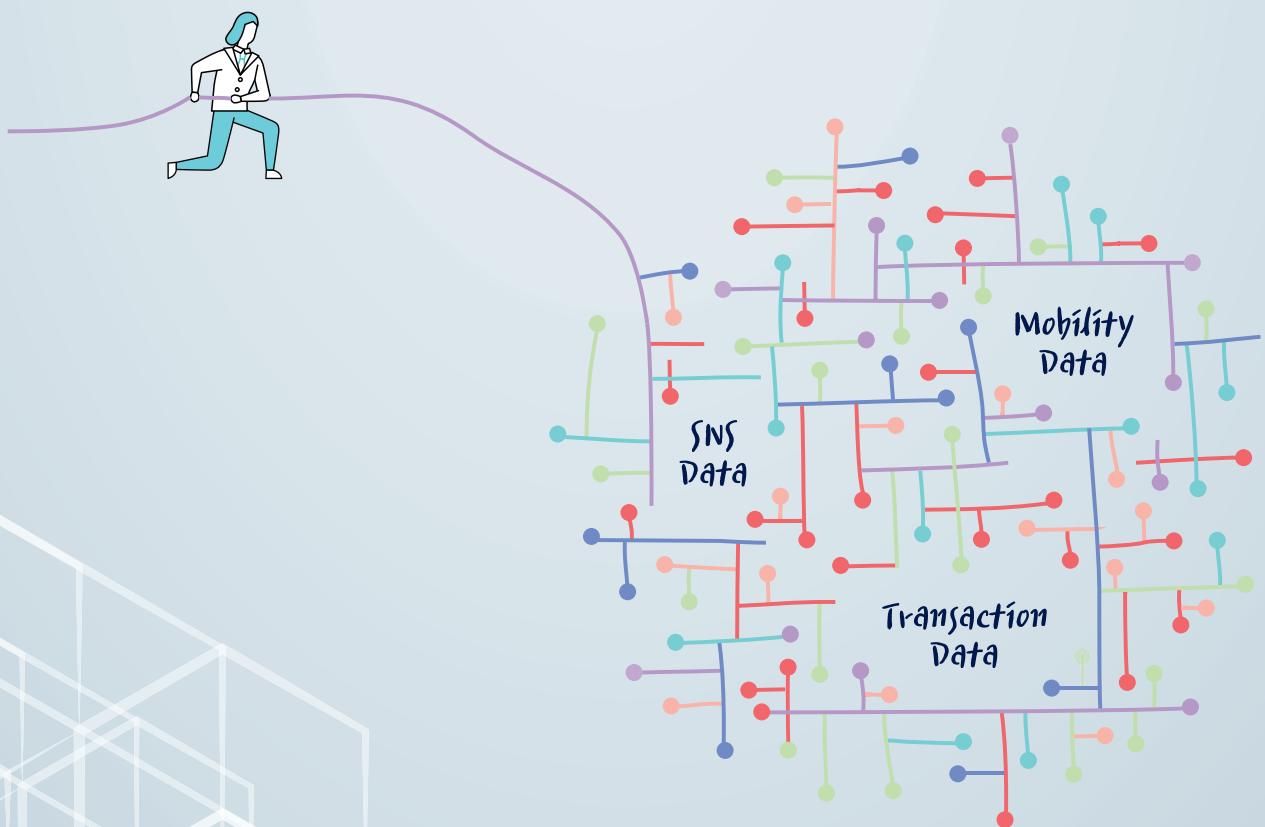
**Thank you.
QnA**

31

▶ 2부: 빅데이터 활용 사례

정치연구와 빅데이터

한규섭 교수, 노선혜 연구원 (서울대 언론정보학과)



빅데이터를 활용한 정치분석

노선혜 (서울대 언론정보학과)

1



빅데이터 여론 추정

2

여론조사심의 위원회 등록 여론조사: 2017년 대선

조사기관	의뢰기관	등록일	구분1	구분2	응답률	문재인	안철수	문재인- 안철수
한국리서치	한국일보	4/10/2017	유무선훈합	전화면접	19.3	37.7	37.0	0.7
리얼미터	이데일리	4/10/2017	유무선훈합	ARS/전화면접	11.8	41.1	34.8	6.3
조원씨앤아이	쿠키뉴스	4/10/2017	유무선훈합	ARS	5.8	40.6	34.4	6.2
알앤써치	(주)데일리안	4/12/2017	무선RDD	ARS	4.2	42.3	37.0	5.3
한국리서치	JTBC	4/12/2017	유무선훈합	전화면접	22.3	38.0	38.3	-0.3
리얼미터	MBN, 매일경제	4/12/2017	유무선훈합	ARS/전화면접	9.8	44.8	36.5	8.3
리서치뷰	프레시안	4/13/2017	무선RDD	ARS	9.8	46	36.5	9.5
한국갤럽	한국갤럽	4/13/2017	유무선훈합	전화면접	23	40	37	3
리얼미터	MBN, 매일경제	4/14/2017	유무선훈합	ARS	10.6	45.4	30.7	14.7

3

역대 총선 출구조사 흐름

	19 th National Assembly ('12)		18 th National Assembly ('08)		17 th National Assembly ('04)	
	S	DP	S	DP	S	DP
Lower Bound	131	131	155	75	92	157
Upper Bound	147	147	178	93	114	182
Election Outcome	152 (X)	127 (X)	153 (X)	81 (O)	121 (X)	152 (X)

4

여론조사는 미국대선을 잘못 예측했는가?

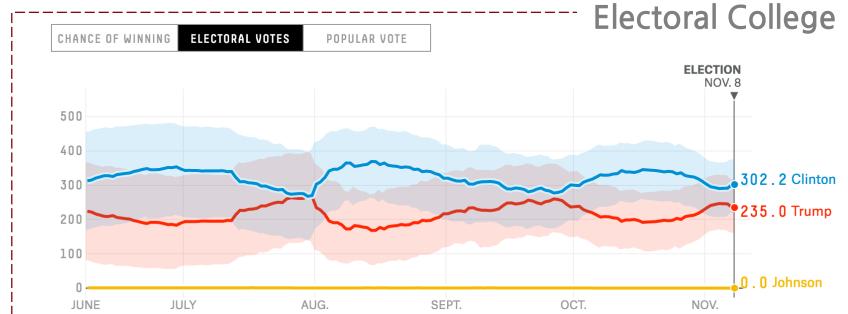
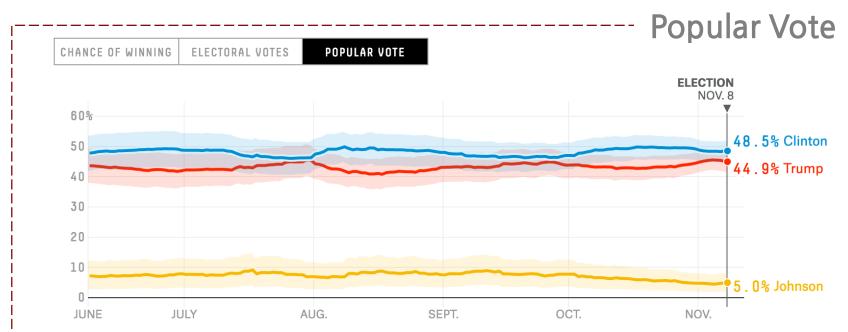
2016 미국 대선 선거인단 예측

			FiveThirtyEight	Wall Street Journal
힐러리	227		302	332
트럼프	304		235	206

발표
3

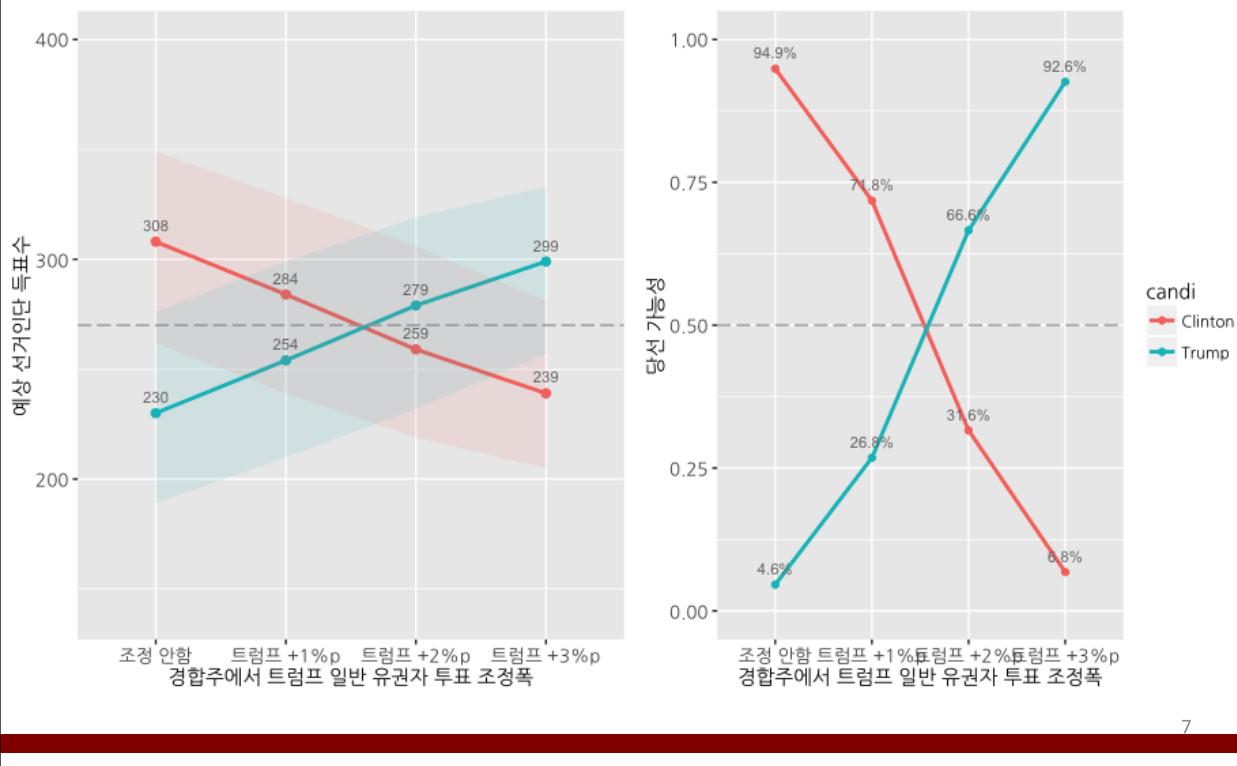
5

여론조사 기반 선거인단 추정 (미국 대선)



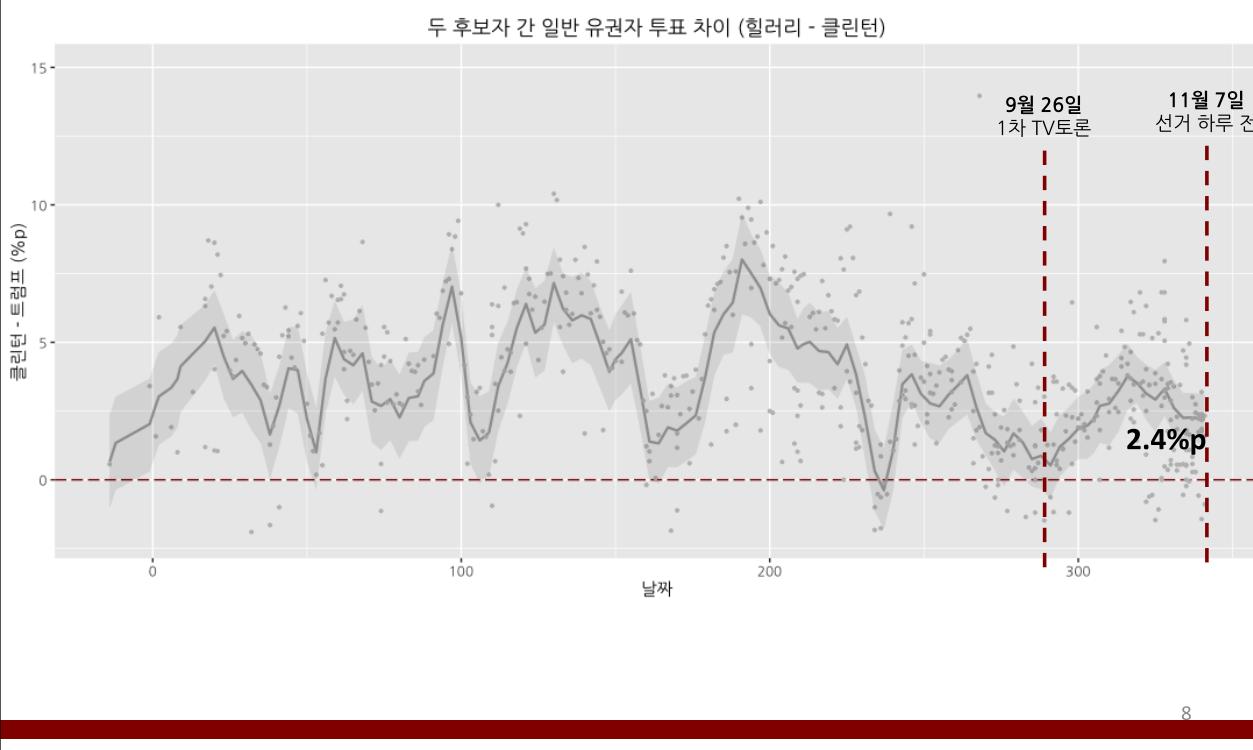
6

2016년 미국 대선 - 선거인단 득표 추정



7

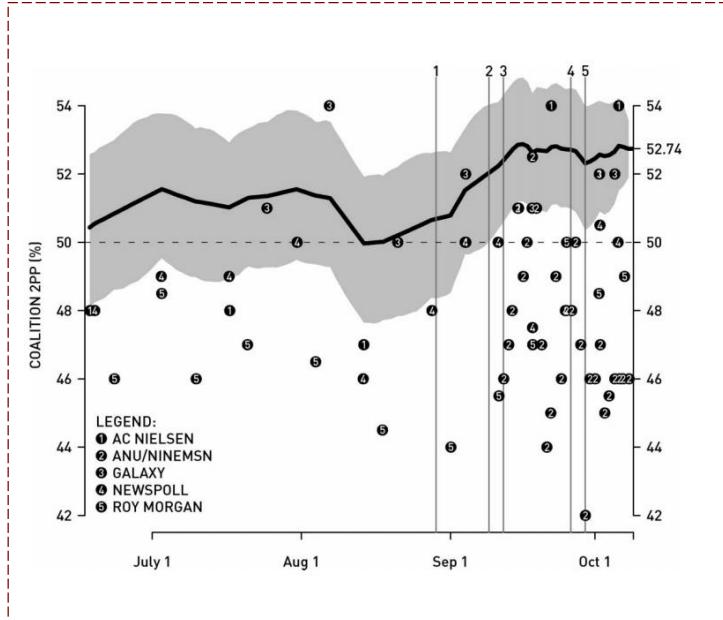
2016년 미국 대선 - 득표율 추정



8

여론조사 어떻게 할 것인가?

Pooling the polls over an election campaigns

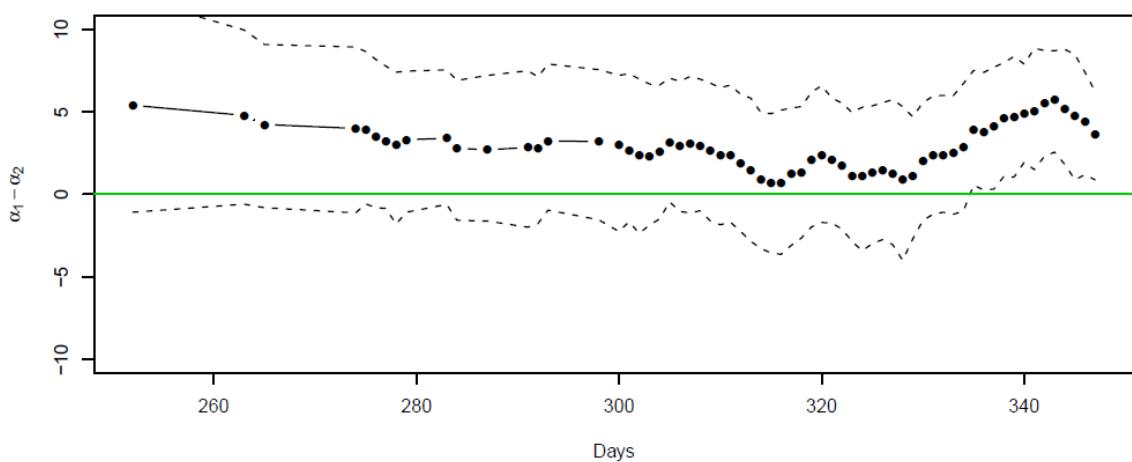


사이먼 잭맨 (Simon Jackman)
호주 국립대학교 정치학과 교수
前 스탠포드대학교 정치학과 교수

9

2012년 대선 예측

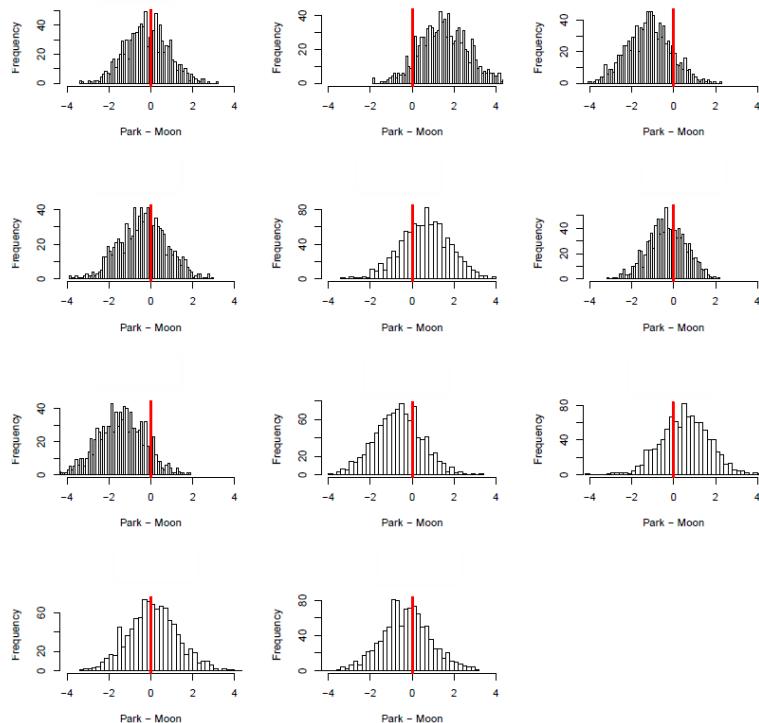
Park - Moon



10



2012년 대선 여론조사 - 조사기관 효과



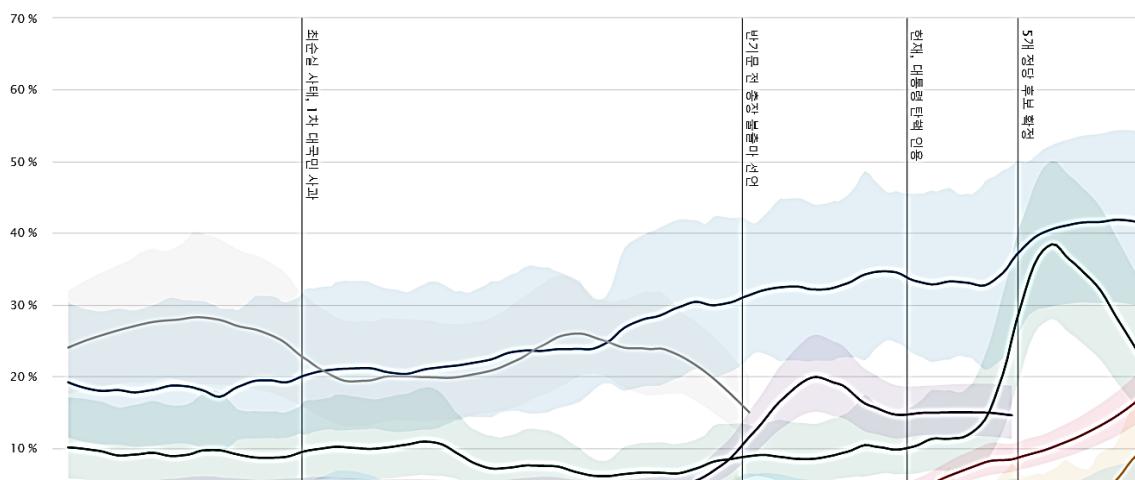
11



pollab 2017 대선 지지율 지수



증합 그래프 지역별 그래프 정당별 그래프 이념별 그래프 연령별 그래프 정당지지 그래프



12

2017 대선 지지율 지수: 최종결과

여론조사 공표 금지 직전 최종

서울대 pollab

	여론조사 공표 금지 직전	등락 (5/2대비)	등락 (4/28 대비)
문재인	41.1%	0.2%p▲	0.9%p▼
안철수	20.8%	0.4%p▼	5.1%p▼
홍준표	18.1%	1.1%p▼	3.4%p▲
심상정	8.6%	1.9%p▼	0.9%p▲
유승민	5.1%	0.1%p▲	0.0%p—

여론조사 공표 금지 직전 최종 vs. 실제 결과

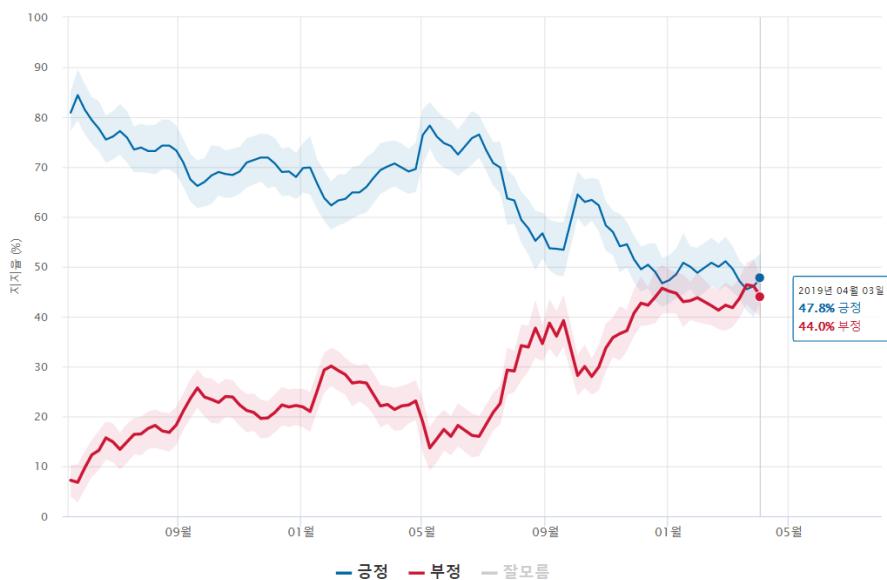
서울대 pollab 지지율 지수

	여론조사 공표 금지 직전	실제결과	오차
문재인	41.1%	41.1%	0.0%p—
홍준표	18.1%	24.0%	5.9%p▲
안철수	20.8%	21.4%	0.6%p▲
유승민	5.1%	6.8%	1.7%p▲
심상정	8.6%	6.2%	2.4%p▼

발표
3

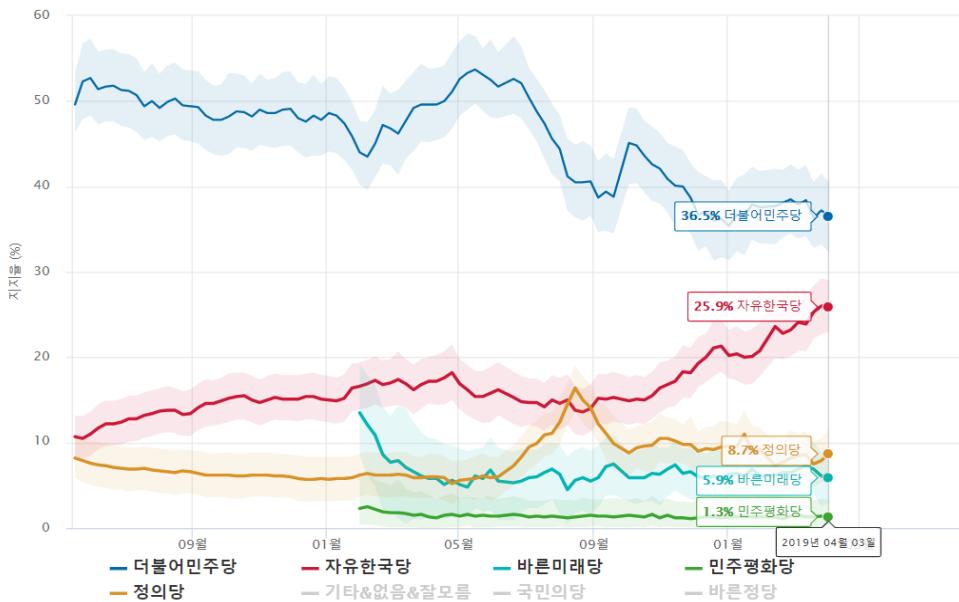
13

대통령 국정운영 평가



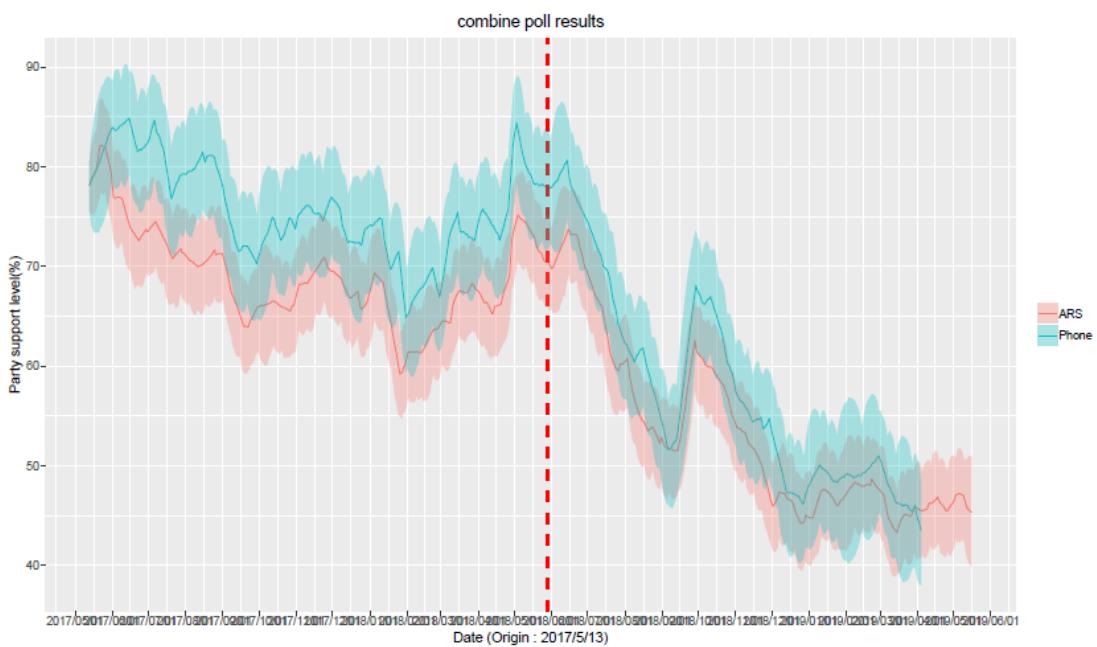
14

정당 지지율 평가



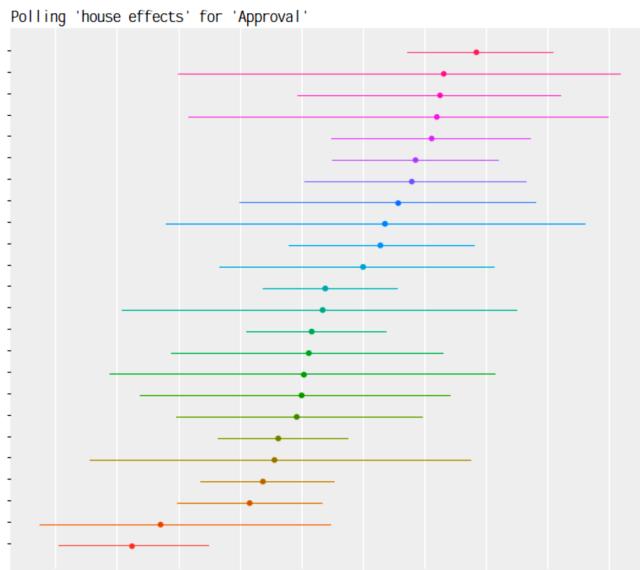
15

대통령 지지율 조사: ARS vs. 면접조사



16

대통령 지지율 조사: 조사기관 효과 (House Effects)



발표
3

17

텍스트 빅데이터 분석

18

우리동네 공약지도 (중앙선거관리위원회): 언론보도 + 지방의회 회의록

언론기사 분석

지방의회 회의록 분석

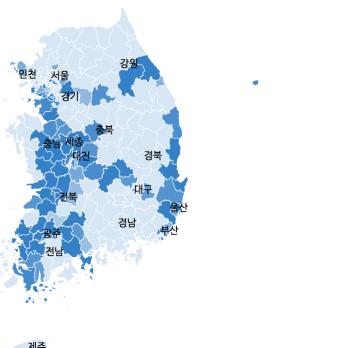
전국

#교육 1위	#소방서 2위	#청소년 3위	#안전 4위	#학교 5위	#경찰 6위	#학생 7위	#조류독감 8위	#농협 9위	#보건소 10위	#아파트 11위
#기업 12위	#장애인 13위	#일자리 14위	#국비 15위	#중소기업 16위	#메르스 17위	#어린이 18위	#전통시장 19위	#중국 20위	#농기센터 21위	
#여성 22위	#이르신 23위	#회계 24위	#구제역 25위	#고등학교 26위	#긴장 27위	#관광객 28위	#강학금 29위	#관광 30위	#환경 31위	
#농업 32위	#봉사활동 33위	#기족 34위	#청년 35위	#국회 36위	#문화재단 37위	#저소득 38위	#학신도시 39위	#학부모 40위	#폭염 41위	
#소년체전 42위	#가을 43위	#안전점검 44위	#노조 45위	#어린이집 46위	#대학생 47위	#도민체전 48위	#아이 49위	#KTX 50위		

시·도별 비교 (KTX)



구·시·군별 비교 (KTX)



19

지역별 이슈 지도

“미세먼지”

구·시·군별 이슈 지도 (미세먼지)



“일자리”

구·시·군별 이슈 지도 (일자리)



“창조경제”

구·시·군별 이슈 지도 (창조경제)

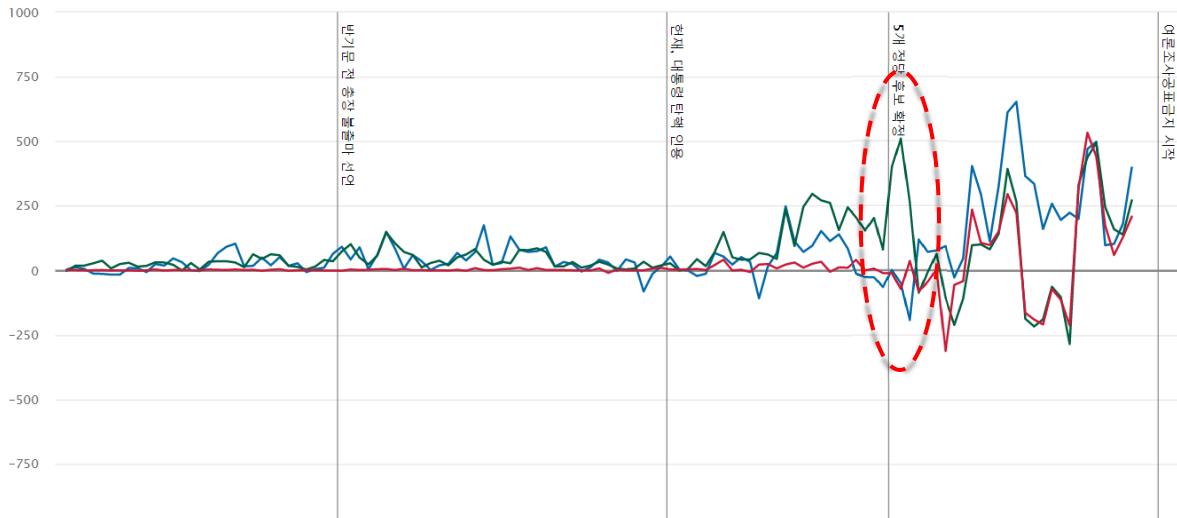


2017 대통령 선거 pollab 보도지수

언론 논조 분석 (기계학습)

pollab

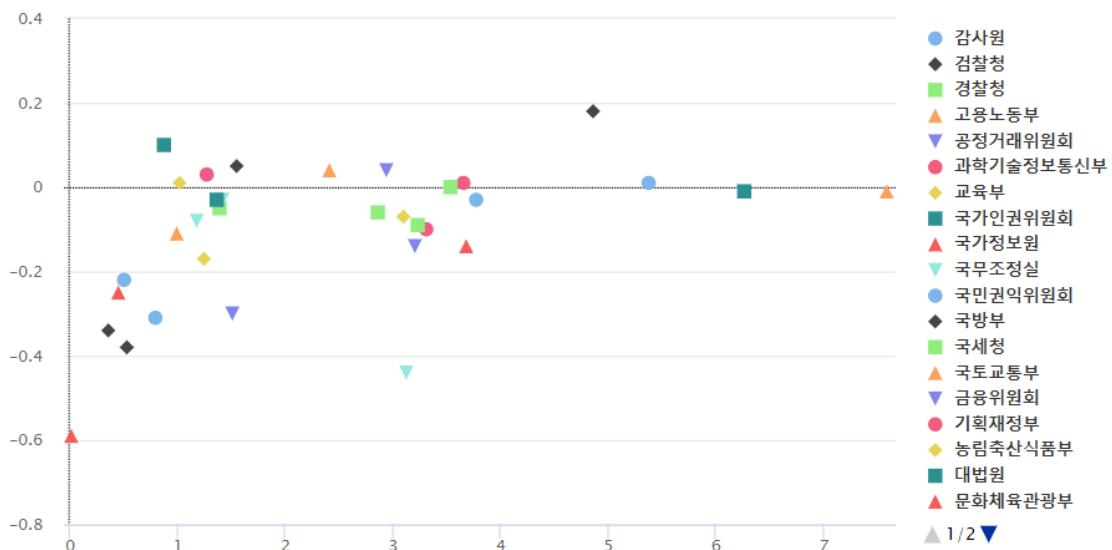
종합보도지수



21

33개 공공기관 신뢰도 조사

공공기관 언론 논조 분석 (기계학습)

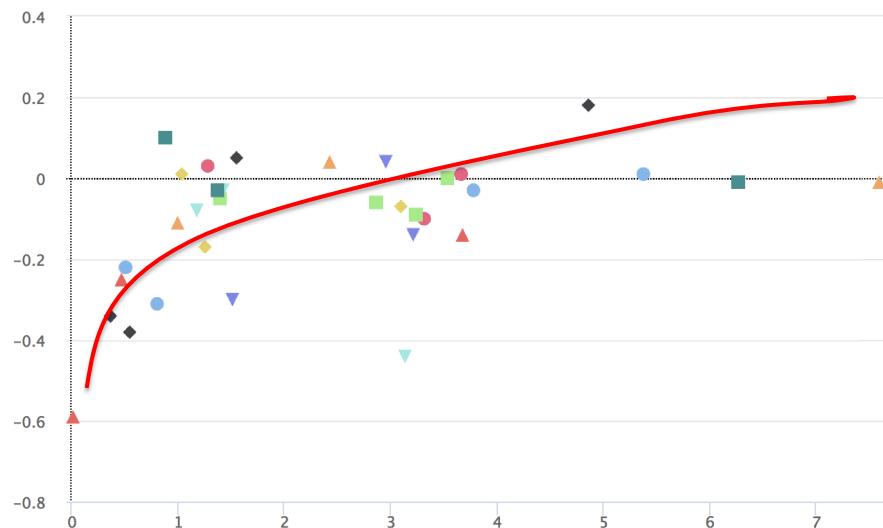


22

33개 공공기관 신뢰도 조사

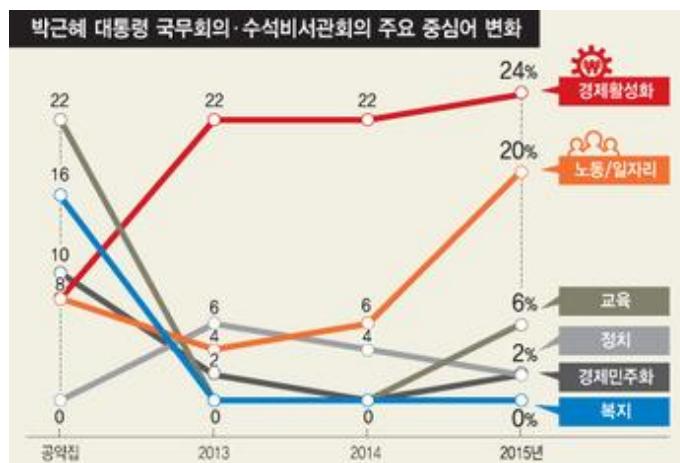
공공기관 언론 논조 분석 (기계학습)

공공기관 언론보도 논조 - 언론논조와 공공기관 평가



23

박근혜 대통령 회의 및 신년사 분석



24

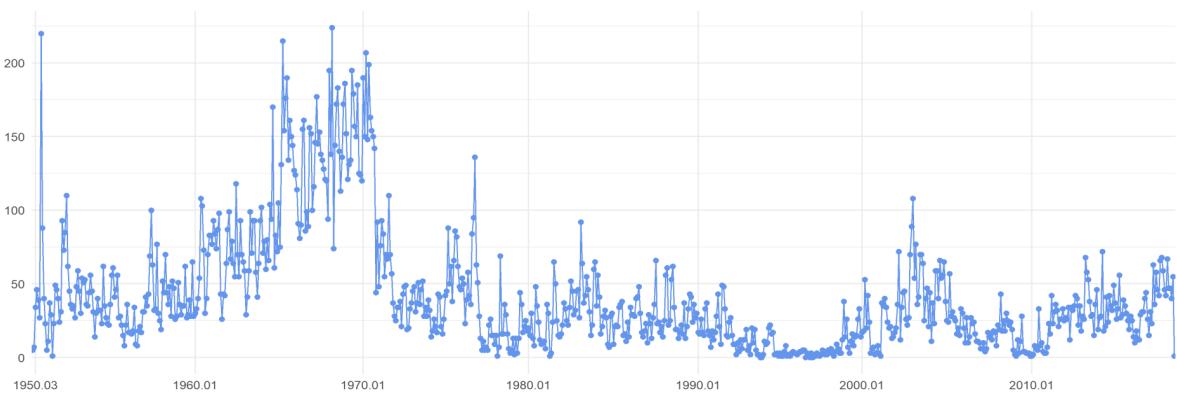
박근혜 대통령 회의 및 신년사 분석

발표
3

25

북한 노동신문 및 조선중앙통신 분석

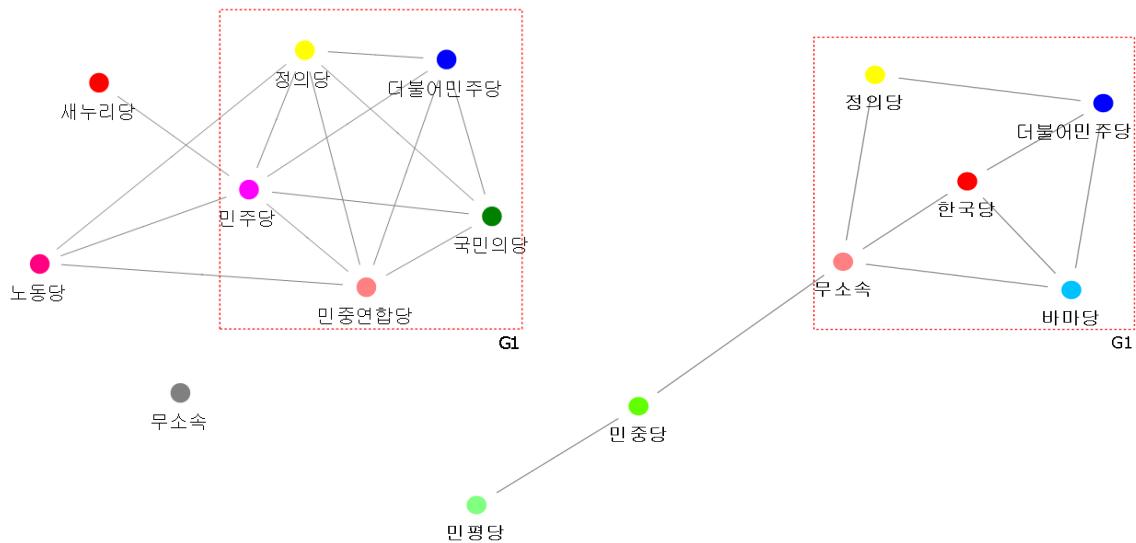
미국에 대한 노동신문 논조



26

국회의원선거, 지방선거 TV 토론 분석

후보자 및 정당 간 네트워크



27

나에게 딱 맞는 후보찾기

28

“나에게 딱 맞는 후보찾기”

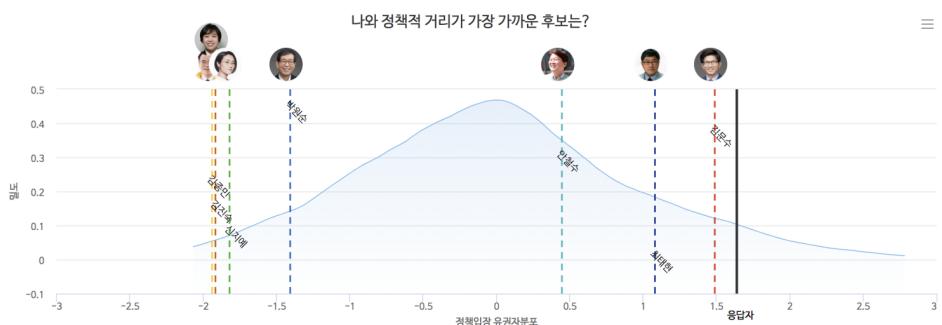
지금까지 여러분의 거주 지역을 선택하면 설문을 시작합니다.



1. 귀하께서는 평균값 중 어떤 줄값이 더 귀하의 일상에 가깝다고 생각하십니까?
 평균과 기업 규제는 금융과 이자율 추세가 일치한다.
 평균과 기업 규제는 노동 투자를 한다.

29

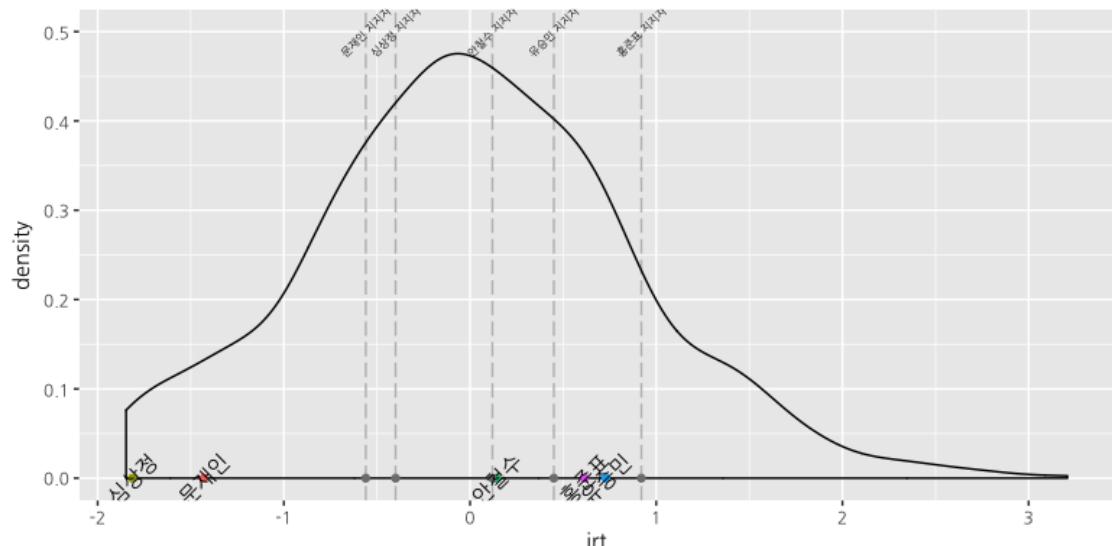
“나에게 딱 맞는 후보찾기”



31개 설문으로 추정한 귀하의 정책입장 점수는 1.64이며,
서울유권자 100명 중 귀하는 95번째로 진보적, 5번째로 보수적입니다.

30

2017년 대선: 정책투표

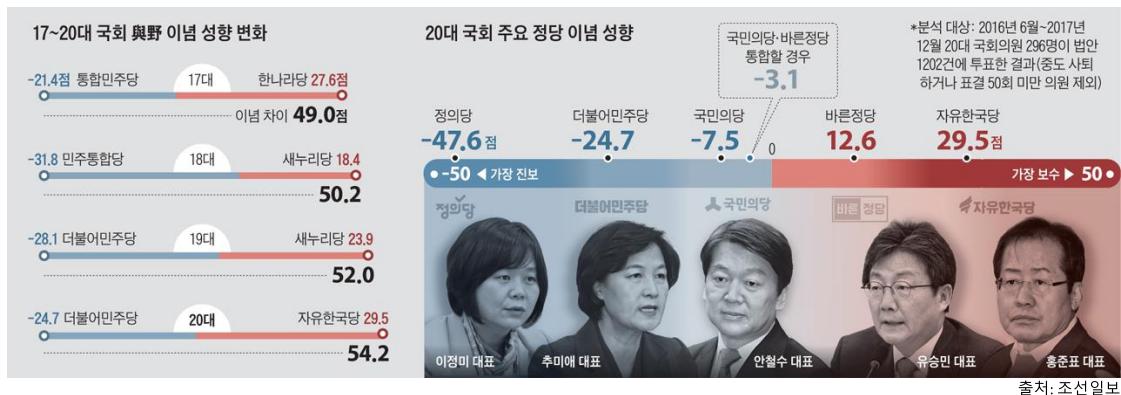


31

국회 표결 및 법안 발의 분석

32

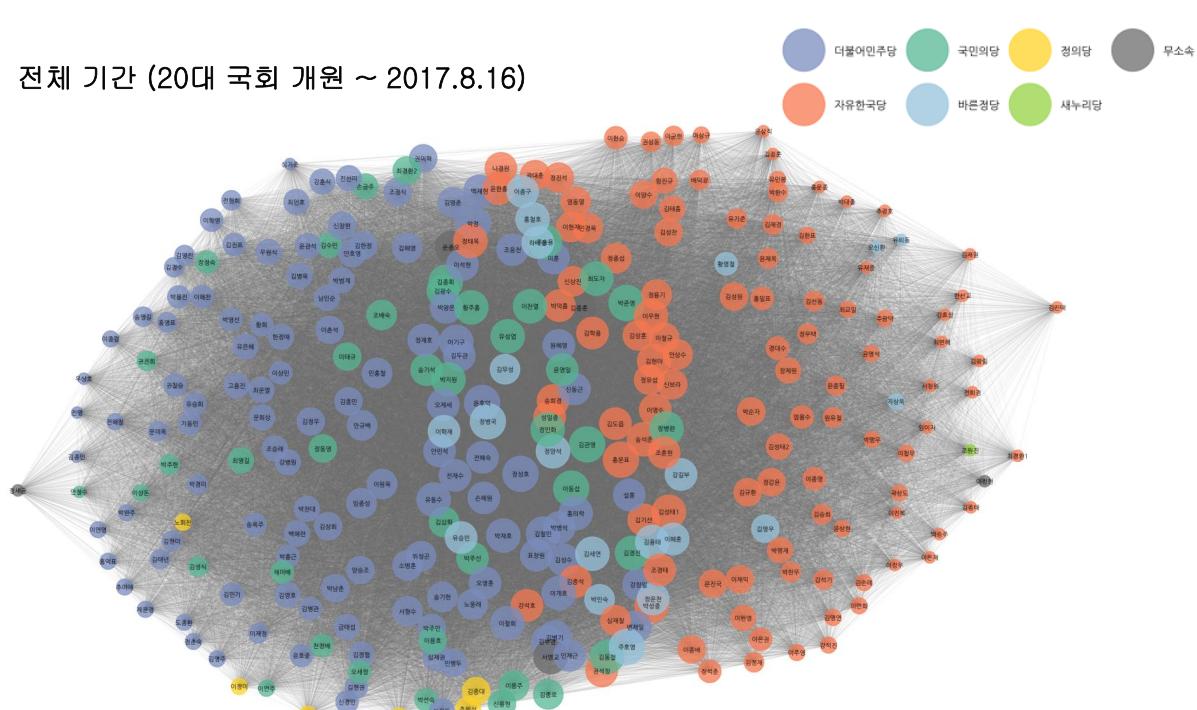
국회 표결 결과 분석



33

공동 발의 분석

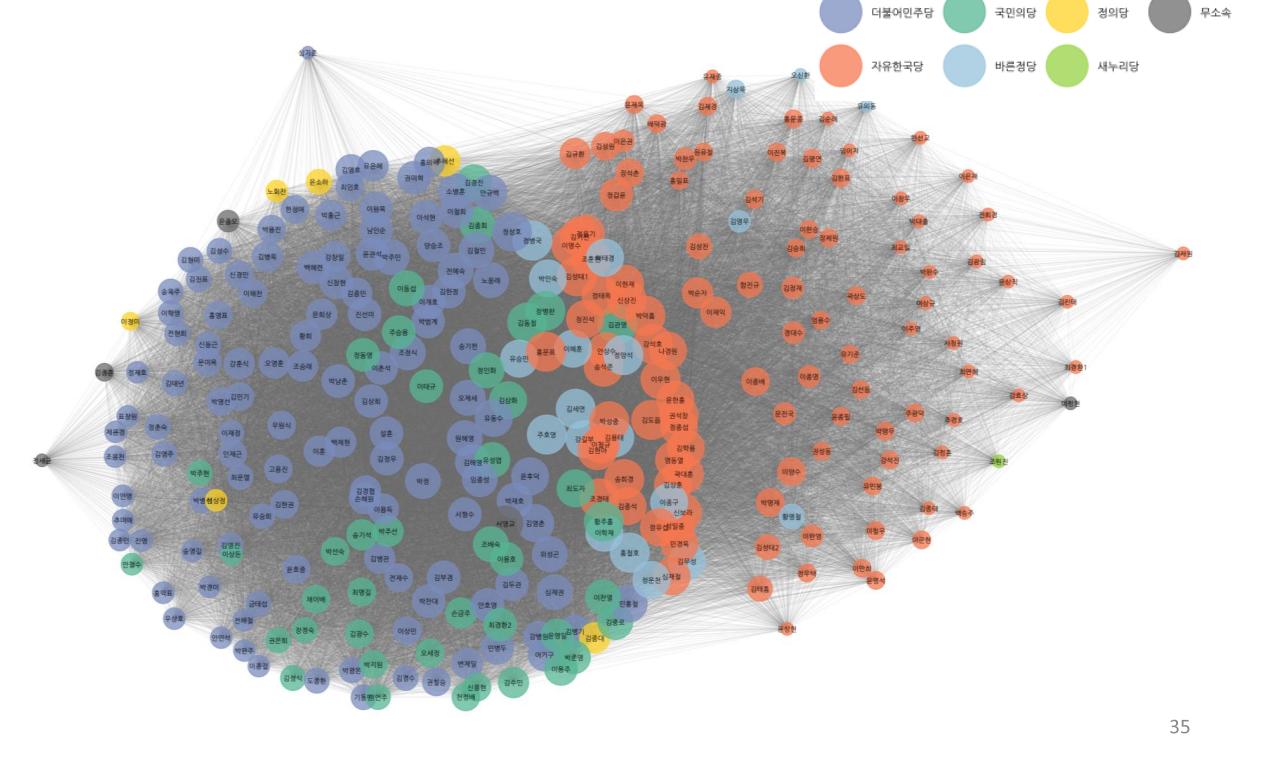
전체 기간 (20대 국회 개원 ~ 2017.8.16)



34

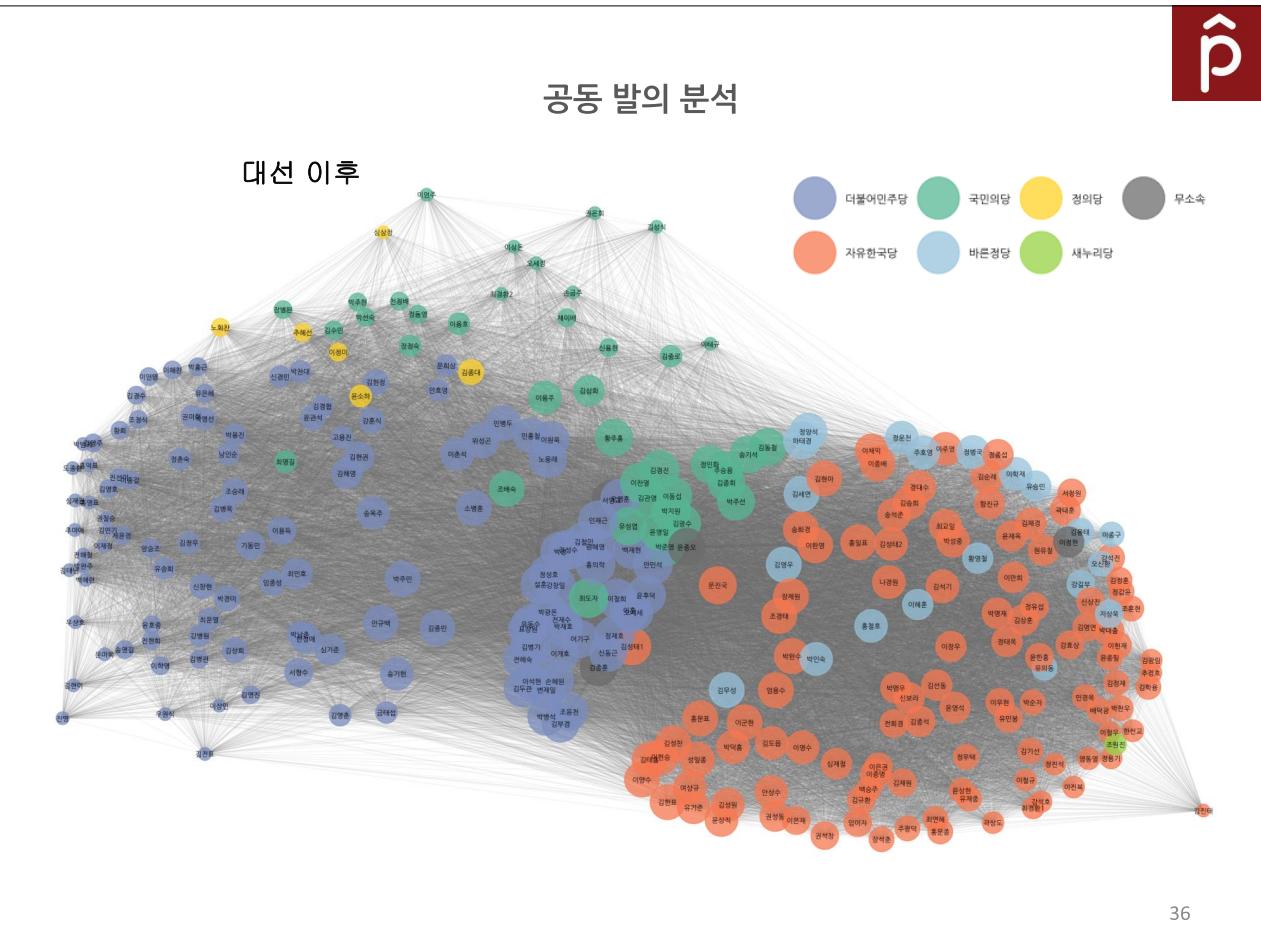
공동 발의 분석

대선 이전



공동 발의 분석

대선 이후

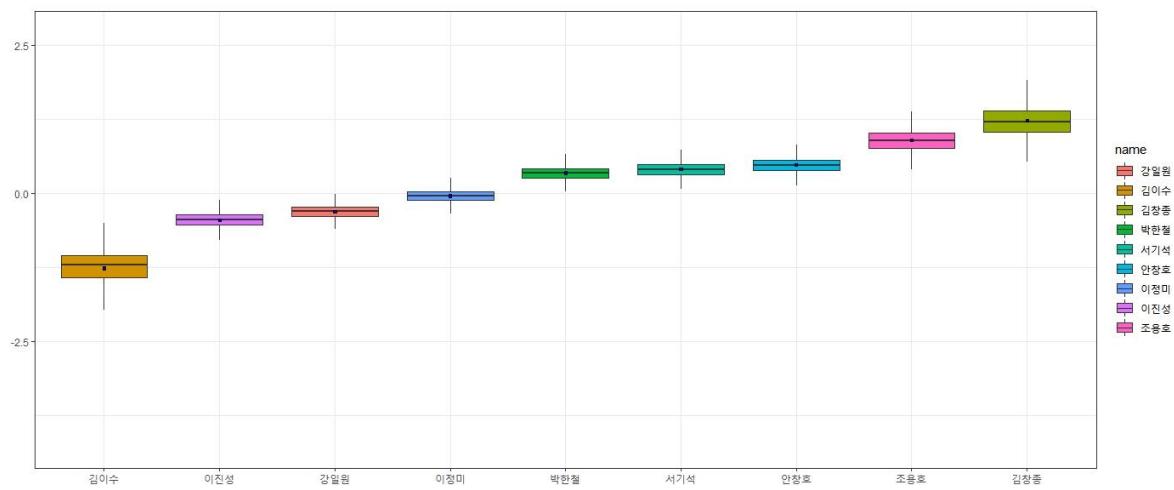


법원 판결 분석

발표
3

37

헌법 재판관 분석



38

헌법 재판관 분석

박근혜 대통령 탄핵 심판 헌법재판관 구성 및 성향



출처: KBS



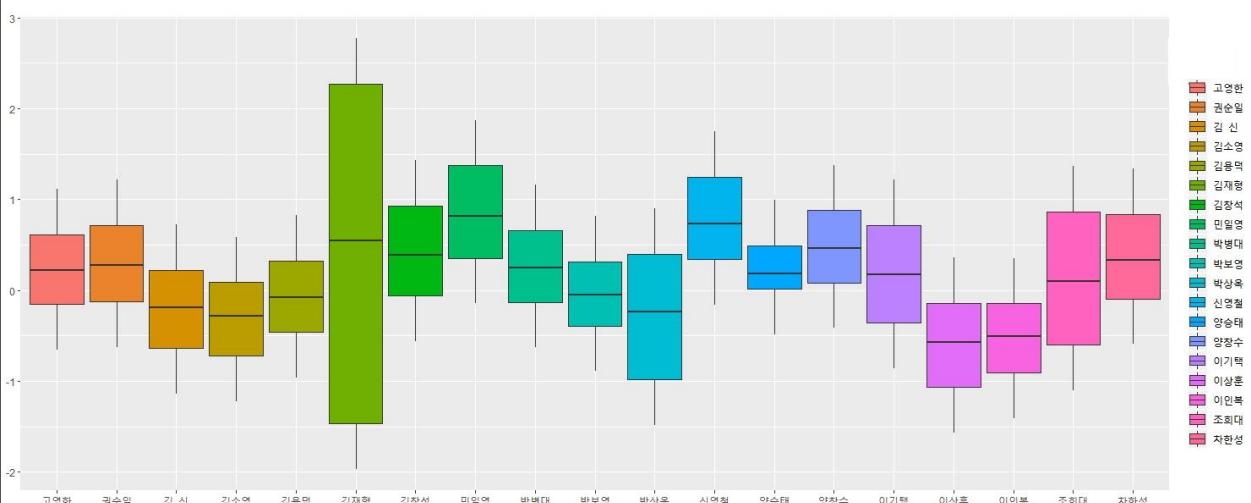
진보적

보수적

39

대법원 전원합의체 판결 분석

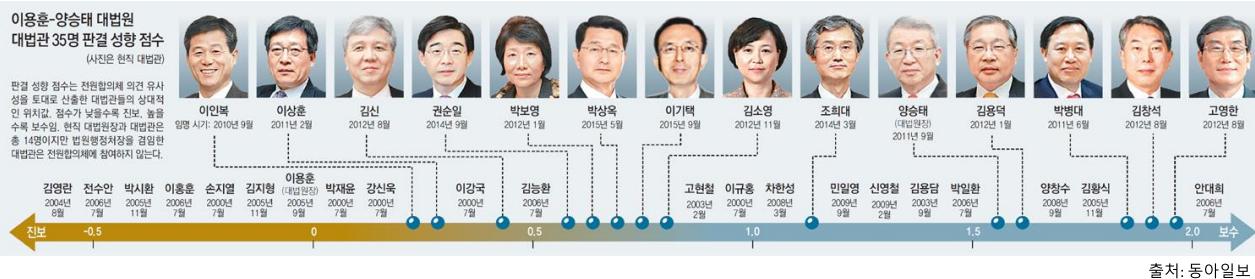
대법관 35명 판결 성향 점수



40

대법원 전원합의체 판결 분석

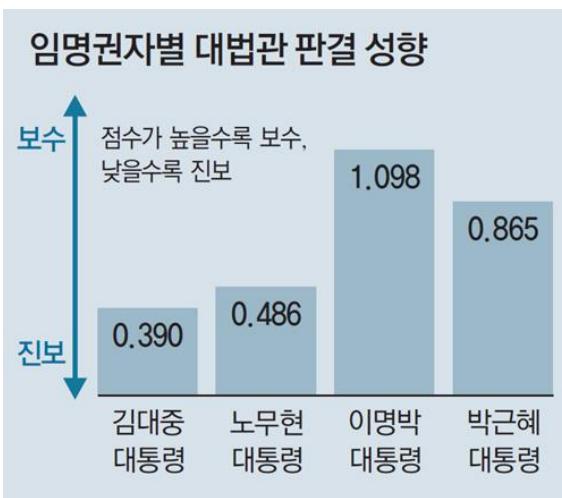
대법관 35명 판결 성향 점수



41

대법원 전원합의체 판결 분석

임명권자별 대법관 판결 성향



42



감사합니다

한규섭 (서울대 언론정보학과)
kyuhahn@snu.ac.kr

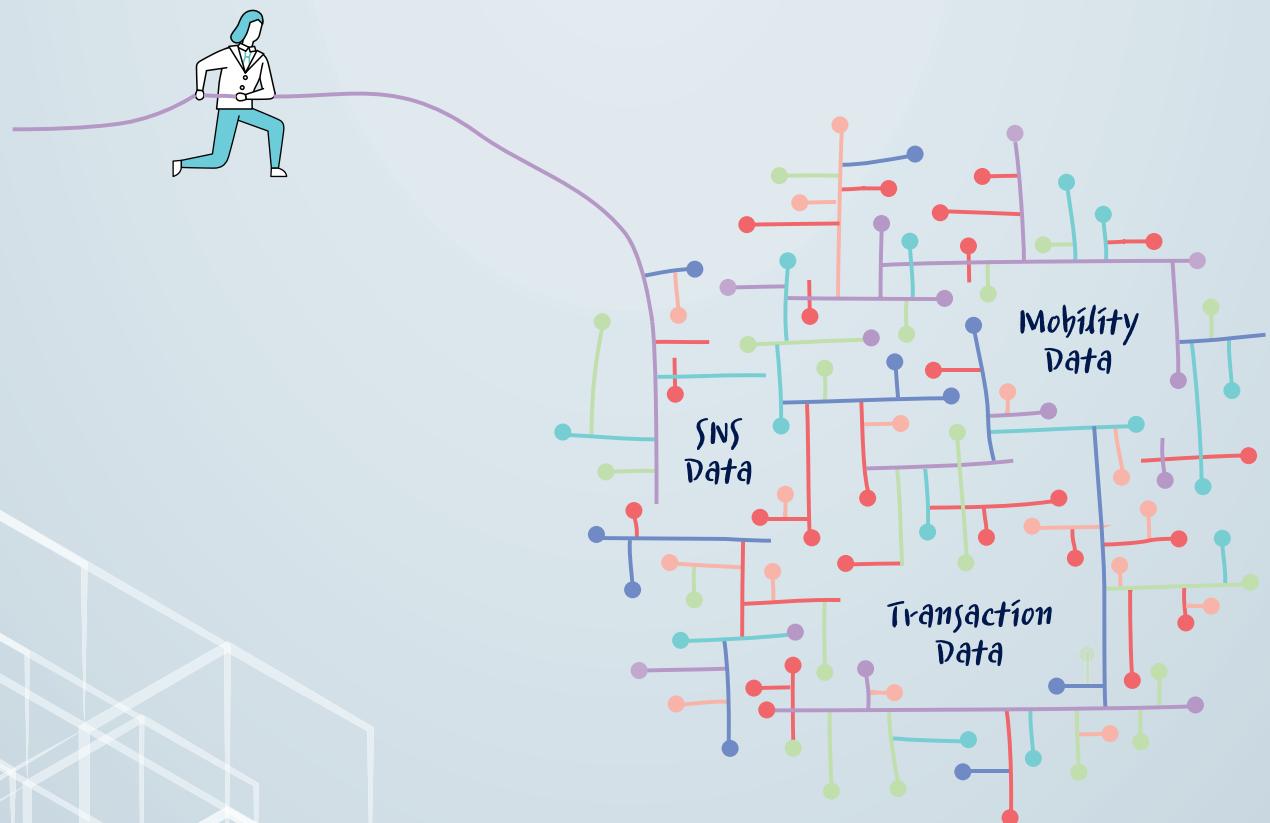
노선혜 (서울대 언론정보학과)
shno205@snu.ac.kr

43

▶ 2부: 빅데이터 활용 사례

문학권력에 대한 세 가지 데이터 분석

이원재 교수 (카이스트 문화기술대학원)



문학권력에 대한 세 가지 데이터 분석

사회, 시장, 문학, 그리고 통계

이원재 KAIST @ KOSSDA 데이터페어 062718

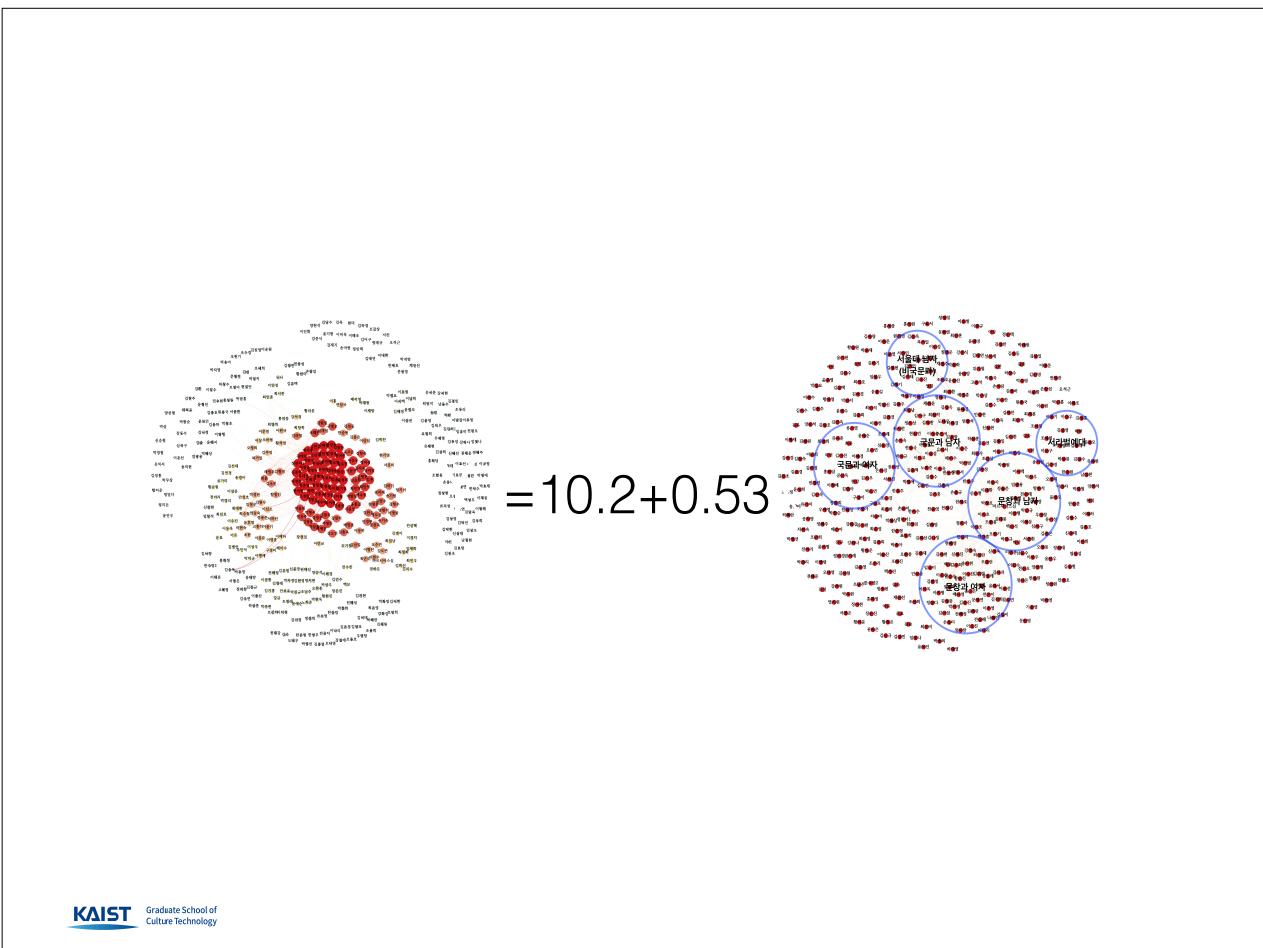
발표
4



사회

(인연 기반 네트워크)





統計



[장정일 칼럼] '지하철 시(詩)'와 문학권력

신경숙 사태 이후 문학권력 논쟁이 분분했지만, 지금까지 논해진 것들은 아무것이나 넣으면 만두 속이 되는 것처럼 편의적이었다. 예컨대 10월 29일자 한국일보 기사(▶ [숫자로 확인된 문학권력의 실체](#))로도 소개된 전봉관·이원재·김병준의 정량적·통계적 연구는 문학권력에 대해 아무 것도 말해 주지 못하는 '통계의 미술'일 뿐이다. 거대 출판사의 상업주의 전략을 문학권력으로 적시하고 나면, 지하철 같



염종선. 2015. "창비를 둘러싼 표절과 문학권력론 성찰." 43(4):66-92.

문학권력 관련 연구논문에 대하여

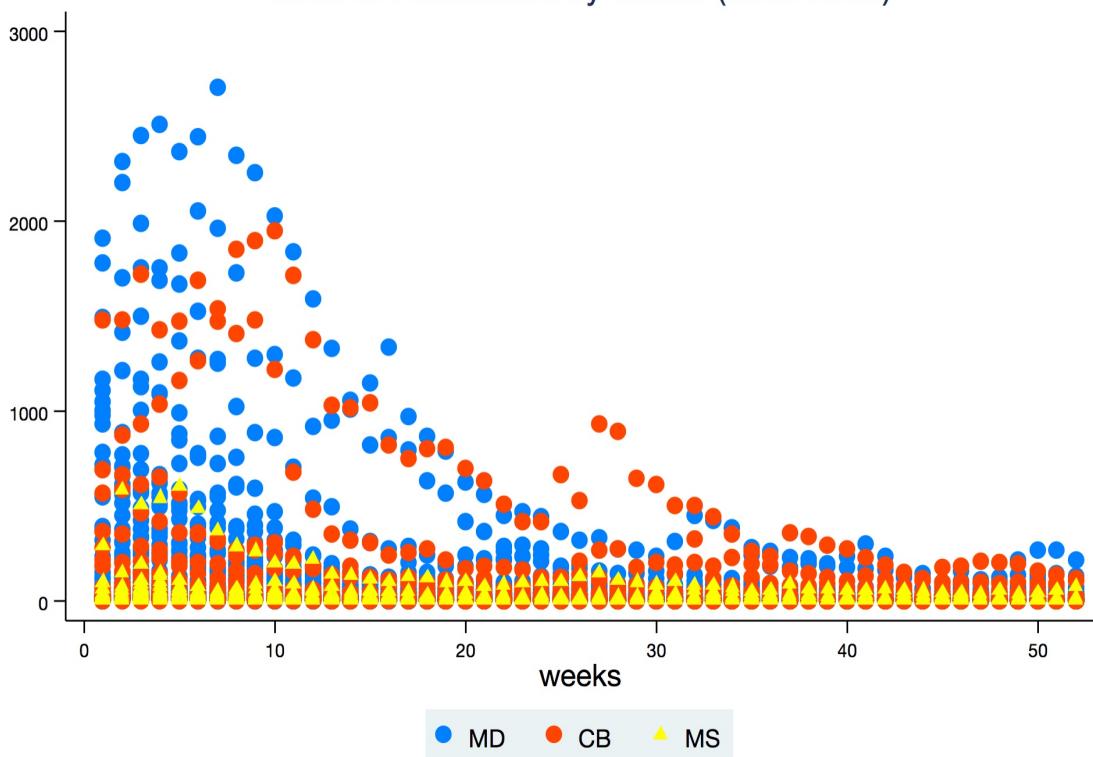
최근 카이스트의 한 연구진은 「문예지를 매개로 한 한국 소설가들의 사회적 지형」이라는 논문을 발표했고 일부 일간지에서 이에 관해 보도한 바 있다.³¹⁾ 1994년부터 2014년까지 21년간 소설을 대상으로

3대 출판사에 관한 내용이 섞여 있어서 일괄적으로 말하기 곤란하지만, 창비에 국한한다면 자사 출신 작가라고 편파적으로 밀어준다는 통념은 사실이 아님이 밝혀졌다. 또한 다른 각도에서 살펴보면, 논문은 “3대 계간지 중복 ‘주도작가’ 총 26명”³²⁾의 명단을 밝히고 이들이 3대 계간지에 작품을 게재하거나 비평적으로 언급되고 단행본 출판을 한 현황을 적고 있는데, 문예지들에서 당대의 주요 작가들에 발표되면 을 주고 비평적으로 언급한 것이 왜 문제인지는 밝히지 않았다. 상업적 동기에서 엉터리 작품들을 밀어주었다면 문제이겠지만 그에 대한 언급이 없는 것은 양적인 분석방법의 한계로 여겨진다. 더군다나 “3대 계간지 모두에 중복된 주도작가가 26명이나 된다는 것은 그만큼 소설 분야만큼은 3대 계간지가 획일화되었음을 의미한다”³³⁾라는 단언은 무슨 뜻인지 알 수 없다. 지난 21년간 얼마만큼의 모집단 수 중에서 26명 정도면 ‘획일화’를 확정할 수 있는 수준이라는 것인가.³⁴⁾

시장

(교보문고 단행본 판매량)

KYOBO Book Sales by Weeks (2010-2015)

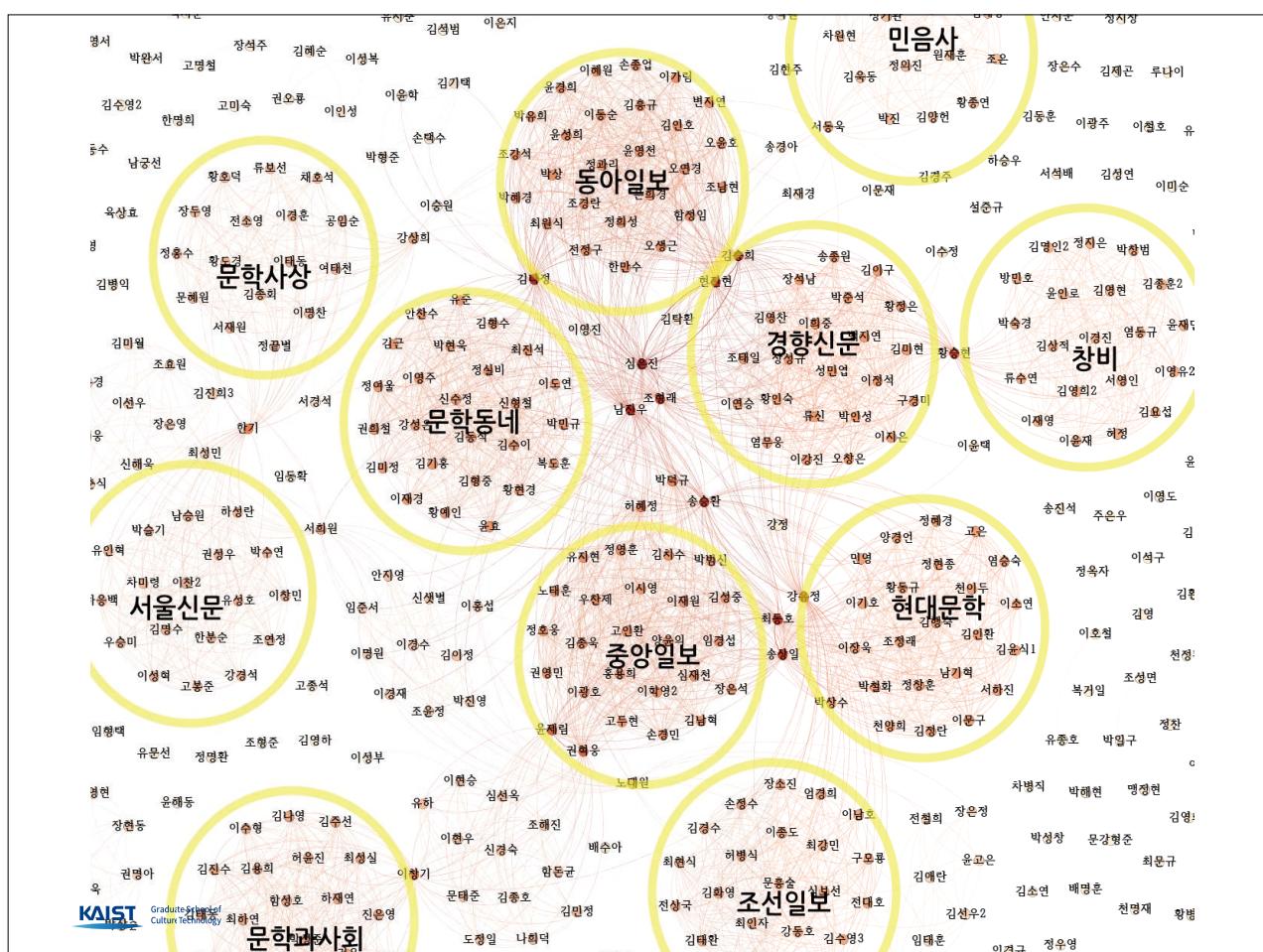
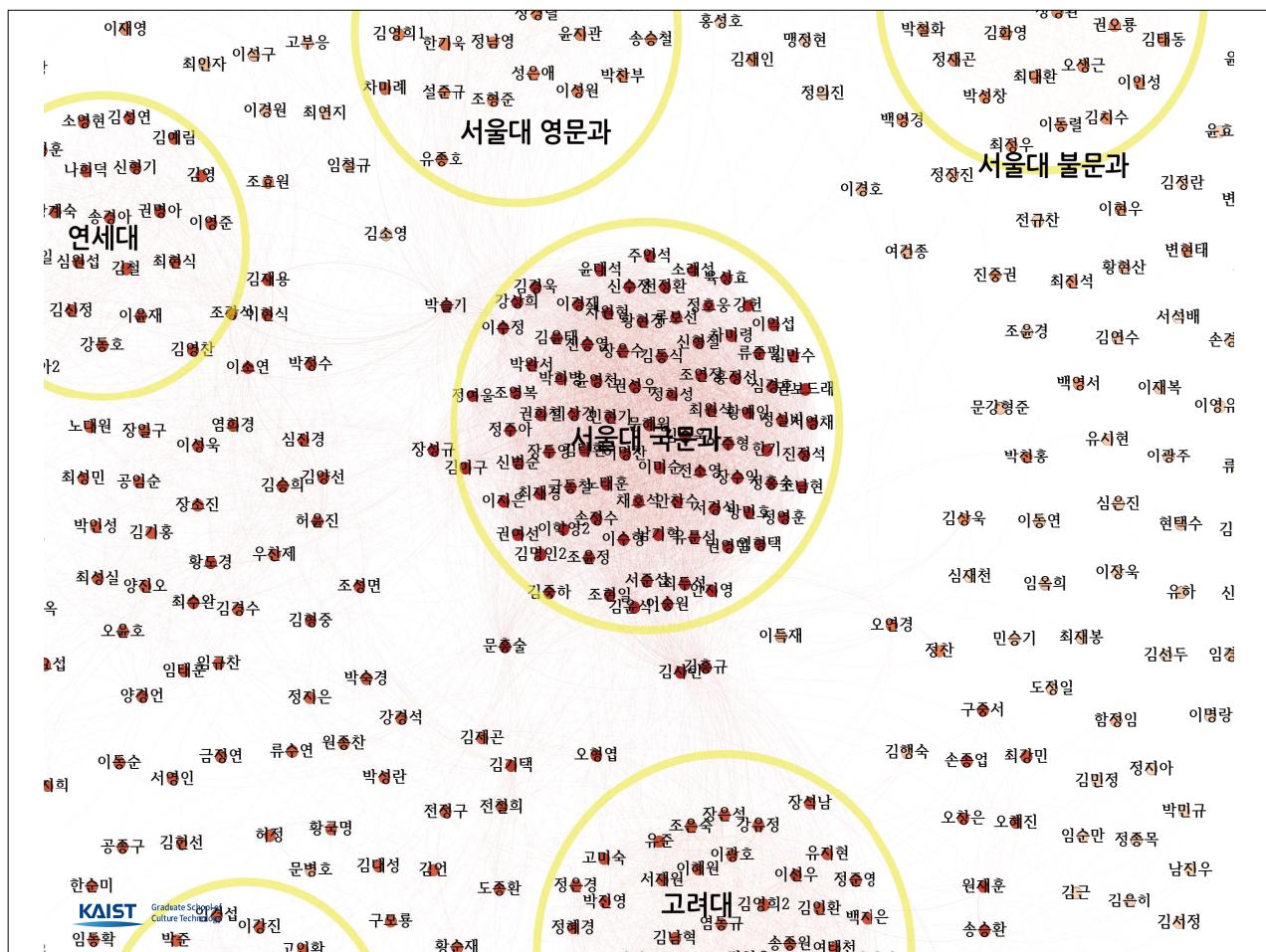


Source: KYOBO BOOK CENTRE 2016

KAIST Graduate School of
Culture Technology

$$Y_{k,w+1} = Y_{k,w} + S_{k,w}\phi + \delta O_{ij} + W_j\beta + \mu_k + \omega_k + \tau_k + \varepsilon_{k,w+1}$$

KAIST Graduate School of
Culture Technology



		(이원재, 김병준, 전봉관 2016)		
		(1)	(2)	(3)
종속 변수: 판매량 증가분 (다음 주(週))				
판매량 증가분 (현재 주)		0.044*** (0.007)	0.042*** (0.007)	0.042*** (0.007)
판매량 (현재 주)		-0.077*** (0.002)	-0.076*** (0.002)	-0.076*** (0.002)
블로그 수 (현재 주)		-0.117 (0.095)	-0.113 (0.095)	-0.122 (0.095)
언론 기사수 (현재 주)		0.753** (0.237)	0.735** (0.237)	0.746** (0.237)
창비 평론수 (현재 주)			-19.509*** (3.030)	-29.243*** (4.803)
문사 평론수 (현재 주)			-2.818 (2.696)	-7.188 (4.294)
문동 평론수 (현재 주)			2.918 (2.115)	-0.356 (3.423)
작가-창비비평가				10.194** (3.924)
인구학적 유사성				4.322 (3.399)
작가-문사비평가				2.952 (2.433)
인구학적 유사성				

문학

(개념어 추출)

From the practical viewpoint, it is immediately apparent that **superhuman intelligence** and intuition are necessary to isolate this entire aggregate of structural relations (the expression of which is essential for understanding the works in question) simply by means of textual studies, no matter how profound and prolonged they may be.

Lucien Goldman 1980

3대 문예지 유사도 그래프_TFIDF: 1995~2015

