

「공식통계 분야에서 BIG DATA」 참가 결과보고

I 출장개요

- 기간: 2014. 10. 27(월) ~ 10. 31(금)
- 장소: 중국 베이징
- 출장목적: 「공식통계에 있어서의 빅데이터(International Conference on Big Data for Official Statistics)」 국제회의 참석
- 출장자(2명): 조사관리국장 강창익 / 복지통계과 전용수 주무관

II 회의개요

- 회의 개요
 - (주 제) 「공식통계에 있어서의 빅데이터」
 - (주 관) 유엔통계처(UNSD), 중국 국가통계국 공동 주관
 - (일시/장소) '14.10.28~30(3일간)/중국 베이징
 - (내 용) 빅데이터와 개인정보보호, 품질유지 방안, 적용기술, 법적 측면 등 발표 또는 토론을 통해 각 국가의 전문지식 또는 경험을 공유
- 회의일정

2014. 10. 28. 화요일		
구분	시간	의제
오프닝 및 패널 토론	09:00- 10:30	<ul style="list-style-type: none"> ▪ 개회사 : Mr. Xie Hongguang (Deputy Commissioner of NBS, China) ▪ 기조연설 <ul style="list-style-type: none"> - 중국 공식통계에 있어 빅데이터 연구 및 적용(중국 통계청 Commissioner) - 공식통계에 있어 빅데이터와 프라이버시 문제(UNSD Director) - 빅데이터와 공간 정보(Greg Scott)
	10:30- 12:00	<ul style="list-style-type: none"> ▪ 패널 토론 : 빅데이터와 국제 통계적 협의체 <ul style="list-style-type: none"> - 참여 : UNSD, UNECE, Eurostat 등
	12:00- 14:00	▪ 점심 휴식
오후 프로그램	14:00- 17:00	<ul style="list-style-type: none"> ▪ 토론 및 프로젝트 워크숍 작성 <ul style="list-style-type: none"> - 주제 : 모바일 기기, GPS 등 장비별 적용 와 위치추적 장치 - 전문가 : Margus Tiru - 사례 <ul style="list-style-type: none"> .모바일 기기와 ICT 지표(ITU) .이동통계를 위한 모바일폰 데이터(ISTAT) .교통통계를 위한 교통데이터(Statistics Netherlands) .국제 거래 통계를 위한 소포 추적데이터(UPU/Global Pulse)

2014. 10. 29. 수요일		
구분	시간	의제
오전 프로그램	09:00- 12:00	토론 및 프로젝트 템플릿 작성 - 주제 : 위성 영상과 다른 공간 정보 - 전문가 : Siu-Ming Tam(ABS) - 사례 .농작물 생산 측정을 위한 위성 영상의 활용(ABS) .농업통계를 위한 위성 영상의 활용(Colombia) .농업통계에서 원격 탐지 기술과 정보(NBS China) .지리정보를 이용한 사회, 경제, 환경통계의 결합(INEGI Mexico)
		12:00- 14:00 점심 휴식
오후 프로그램	14:00- 17:00	토론 및 프로젝트 템플릿 작성 - 주제 : 트위터 및 기타 소셜미디어 - 전문가 : Prof. Andrew Schwartz (University of Pennsylvania) - 사례 .소비자 전망 통계를 위한 트위터 데이터(Stactitics Netherlands) .음식 가격 변동 측정을 위한 트위터 데이터(Global Pulse) .노동통계를 위한 웹 스크래핑(ISTAT) .네트워크 데이터 검색 기반하에서 부동산가격 예측(NBS China) .공식통계를 위한 소셜미디어 데이터의 활용(INEGI Mexico) .Baidu 검색지표의 활용(중국 회사)
		토론 및 프로젝트 템플릿 작성 - 주제 : 거래 및 스캐너 데이터 - 전문가 : Peter Struijs (Stactitics Netherlands) - 사례 .가격통계를 위한 스캐너 데이터(Stactitics Netherlands) .온라인 가격 데이터 캡처(NBS China)

2014. 10. 30. 목요일		
구분	시간	의제
오전 프로그램	09:00- 12:00	발표 (운영자 : Steven Landefeld) - 주제 : 빅데이터 원천의 이익과 도전에서의 공통점 - 내용 .공통 방법적 이슈 및 품질관련 .법적 및 윤리적측면을 고려한 프라이버시 문제 .데이터 획득 및 책임 분배 등 파트너십에 대한 이슈 .IT 이슈 등
		발표 (운영자 : UNSD) - 주제 : 개도국 환경에서 혁신의 소개 - 내용 : 데이터 프로젝트/ 데이터 철학/ 빅데이터 경험 등
	12:00- 14:00	점심 휴식
오후 프로그램	14:00- 17:00	패널 토론 : 빅데이터 관련 향후 계획 - 참여 : UNSD, ABS, Stactitics Netherlands, NBS China등
		폐회사 : UNSD, NBS China

Ⅲ 주요의제 및 회의내용

1 외국의 빅 데이터(BIG DATA) 활용현황

□ 빅 데이터 워크숍에서 제기된 외국의 빅 데이터 활용사례를 3가지 유형*)으로 구분하여 제시

*) (유형 1) 인터넷상의 온라인 거래 자료, (유형 2) 개인 간 통신과 관계를 나타내는 모든 자료 (유형 3) 기타, 사회구조기능의 흐름을 나타내는 자료

Big data 원천	활용분야	주요내용	
유형 1 : 온 라인 거래자료	구글, Baidu 등 인터넷 정보	<ul style="list-style-type: none"> • 노동통계 • 물가통계 작성 <ul style="list-style-type: none"> • 구글 트렌드를 활용한 노동시장 예측 조사(이탈리아) • web-crawler 기술을 이용한 소비자물가지수(중국), Baidu 서칭 자료를 활용한 집값예측 (중국) 	
유형 2 : 개인 간 통신 과 관 계 를 나타내는 자료 및 위성자료	통신자료 (Mobile positioning data)	<ul style="list-style-type: none"> • 관광통계 • 교통량통계 • 인구통계 <ul style="list-style-type: none"> • 외국인의 로밍정보를 활용하여 관광객 선호지역, 이동경로 파악(UN 연구) • 국내휴대폰 이용자의 이동경로를 활용하여 교통 관련 통계작성(이탈리아) • 특정지역의 휴대폰 위치정보를 실시간으로 활용하여 실시간 인구파악 (네델란드) 	
	Social Media (페이스북, 트위터)	<ul style="list-style-type: none"> • 보건통계 • 의식조사 	<ul style="list-style-type: none"> • 소셜미디어 정보(메시지, 사진 등)를 활용하여 HIV, Flu 만연 관련 통계, 심장병 사망 가능성 통계 작성(미국) • 소셜미디어 정보를 이용한 소비자심리지수작성(네델란드), 트위터 언어를 이용한 삶의 질 만족도 관련 통계 작성 등(멕시코)
	위성정보 (Satellite Imagery)	<ul style="list-style-type: none"> • 농업통계 • 기타, 총조사 및 환경통계 	<ul style="list-style-type: none"> • 위성사진을 활용한 경작면적 및 작물 생산량 통계(호주, 중국) • GIS를 활용한 경제총조사, 환경통계 (멕시코)
유형 3 : 기타, 사회구조, 기 능 을 나타내는 자료	교통영상 (Traffic loops)	<ul style="list-style-type: none"> • 인구이동통계 • 교통량 통계 <ul style="list-style-type: none"> • CCTV 등 영상정보를 활용하여 인구이동 통계 작성(네델란드) • 고속도로 통행 자료를 활용한 교통량 통계 작성(네델란드) 	

2 빅 데이터(BIG DATA) 활용의 장단점 및 고려사항

□ 빅 데이터 활용의 장점

- (통계생산 효율성 제고) 보다 시의성 있고, 세부적인 데이터를 응답자 부담 없이 낮은 비용으로 획득하여, 맞춤형 자료생산 가능
- (이용자 중심의 통계 제공) 정부 정책결정 및 민간분야 의사결정에 필요한 다양하고 세부적인 통계를 실시간(Real time)으로 제공 가능

□ 빅 데이터 활용의 단점

- (자료 획득의 어려움) 빅 데이터는 주로 사적영역, 특히 이동통신 회사, 구글 등 글로벌 기업이 보유하고 있어 수집하기 어려움
- (구조화 및 분류하기 어려운 자료) 주로 비통계적 목적으로 수집된 자료이므로 구조화, 분류하는데 한계
- (자료의 대표성 문제발생) 수집된 자료가 모집단 전체를 대표하기 곤란한 경우가 많고, 특히, 모바일 통신자료의 경우 이동통신망이 있는 지역에만 적용되기 때문에 대표성 문제가 발생

□ 빅 데이터 도입 시 고려사항

① 방법론적 (Methodological) 고려

- (자료의 대표성 확보) 통계적 방법으로 추출된 표본이 아니어서, 자료의 대표성에 한계(Sample bias)가 발생 ⇒ 모수 추정방법 검토

[예시] 휴대폰 위치정보를 이용한 실시간 인구이동통계 작성 시 휴대폰 미사용자인 유아 및 노인은 제외되고, 청소년, 중장년층이 과대 추출되어 표본의 대표성 (Representative) 저하

○ (Correlation vs Causation 고려) 빅 데이터를 활용하여 작성된 통계와 기존통계 간 상관관계 분석을 통해 활용가능성 검토

- 그러나 상관관계가 높다고 해서 기존의 공식통계를 대체하는 것은 신중해야 함 ⇒ 상관관계뿐만 아니라 현상에 대한 인과관계도 고려

[예시] 트위터, 페이스 북 등 소셜미디어를 활용하여 경기에 대한 소비자심리지수(Consumer Sentiment Index)를 작성하여 기존 통계와 상관관계가 높은 방법론(모델링)을 선택

② 개인정보 보호(Privacy)

○ 소극적동의(Passive consent) vs 적극적동의(Active consent) 필요

- 빅데이터 및 공식통계 특성을 검토하여 정보주체의 동의 필요여부, 소극적동의¹⁾가 필요한지 또는 적극적 동의²⁾가 필요한지를 파악해야함

1) 소극적 동의: 정보주체(개인)와 빅데이터 제공자(기업 등)간 개인정보 활용에 대한 동의를 빅데이터 활용 통계생산에 대한 동의로 간주하는 것

2) 적극적 동의: 정보주체와 통계 생산자(통계청 등)간 개인정보 활용에 대한 별도의 동의가 필요한 것(일종의 참여)

③ 파트너십 구축(Partnership)

○ 빅 데이터 수집 및 처리, 통계 추정방법 등 절차의 투명성 확보와 통계작성기관과 정보제공자(빅 데이터 제공자) 간 긴밀한 협력이 필요

○ 비밀보호 및 신뢰성에 관한 명확하고 강력한 규칙 수립

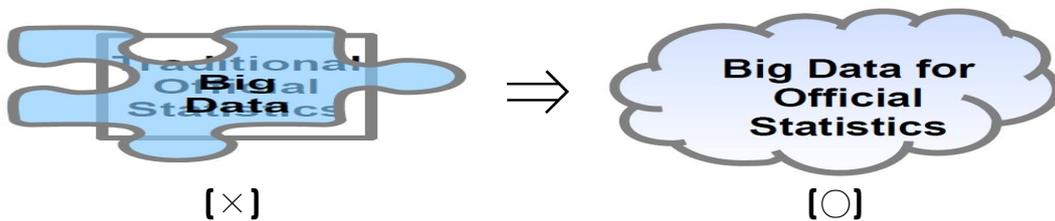
④ 수집된 데이터 관리 및 처리시스템의 구축

○ 대용량 데이터의 빠른 처리를 위해서는 데이터 저장 공간(data storage) 및 분석 인프라(analysis infrastructure) 구축이 요구됨

IV 시사점

□ 빅 데이터 활용을 위한 연구와 인식의 전환 필요

- 전 세계적으로 빅 데이터 활용은 거대한 흐름(Big trend), 우리나라도 공식통계에서 빅 데이터를 활용할 수 있는지 연구할 필요
 - 공식통계의 체계를 유지하면서 빅데이터를 활용하기 보다는 빅데이터의 비구조적인 특성을 고려하여 공식통계에 활용하는 방안 검토 등 유연한 사고로 전환 필요



- 빅 데이터를 공식통계에 활용하기 위해서는, 빅 데이터 활용가능 분야, 활용방법, 자료 수집 및 개인정보 보호 및 빅 데이터 처리 방안 등에 관한 기본전략 수립이 필요

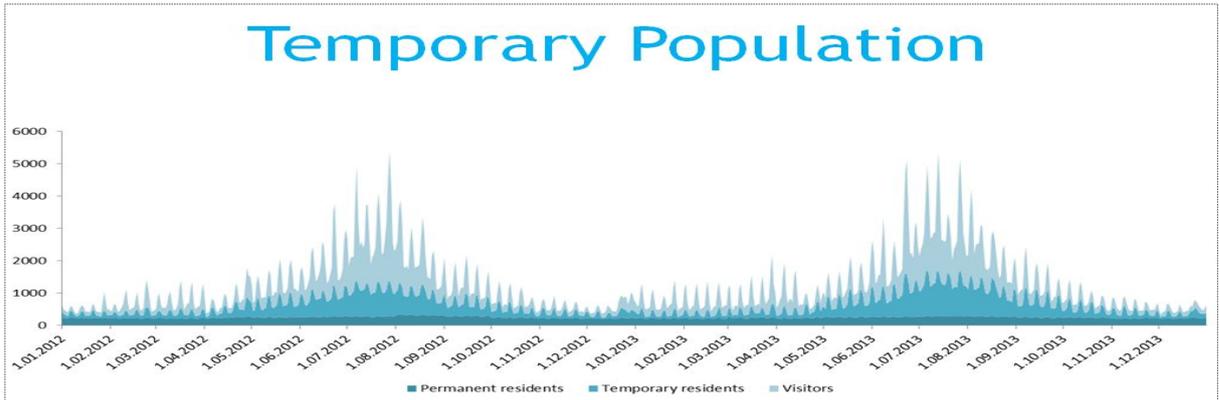
□ 구체적인 사례 제시 - 모바일 정보를 활용한 통계작성

- (모바일 정보의 정의 및 특성) 시공간에서 모바일 장치의 위치를 추적하는 적극적 또는 소극적 위치 데이터

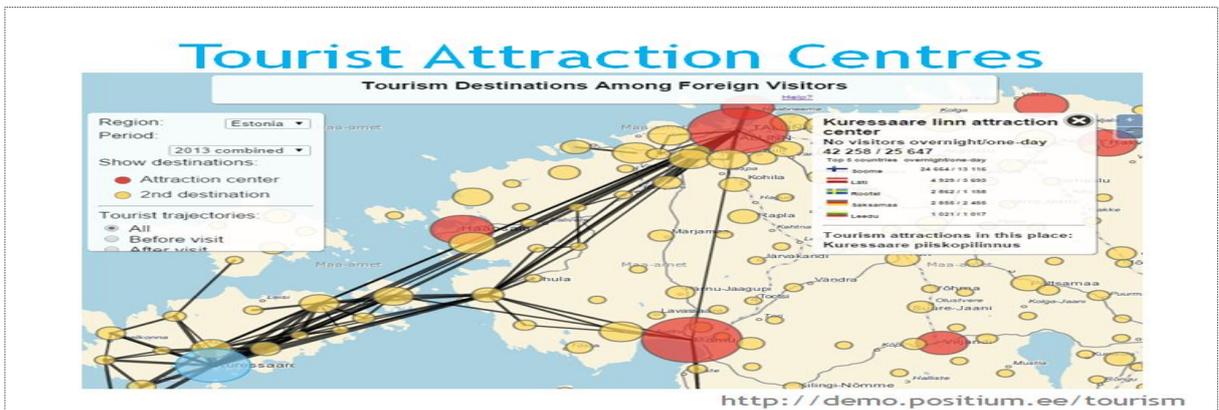
- * 적극적 위치정보(Active positioning) : 휴대폰의 위치정보를 실시간으로 수집 ⇒ 소유자 동의 필요
- * 소극적 위치정보(Passive positioning) : 이동통신회사로부터 받은 휴대폰 위치 정보 ⇒ 소유자 동의 불필요

- 주요 적용 분야

- (인구통계) 휴대폰 보유자의 시간별 위치정보를 활용하여 지역의 실시간 인구통계(daytime population)를 생산하여 도시설계 정책에 활용



- (여행관련 통계) 외국인의 로밍정보를 활용하여 여행객들의 이동 경로 및 체류시간 등의 정보를 분석하여 여행관련 통계 생산



- (교통통계) 국내 휴대폰 보유자의 이동정보(출발지~종착지)를 활용하여 출퇴근 패턴을 분석하고, 교통량 및 주요이동 경로 등 교통통계 생산에 활용

