



# How Big Data Will Revolutionize Government Policymaking?

**Setia Pramana**

Politeknik Statistika STIS

Analisis Big Data dengan R | Pusdiklat BPS

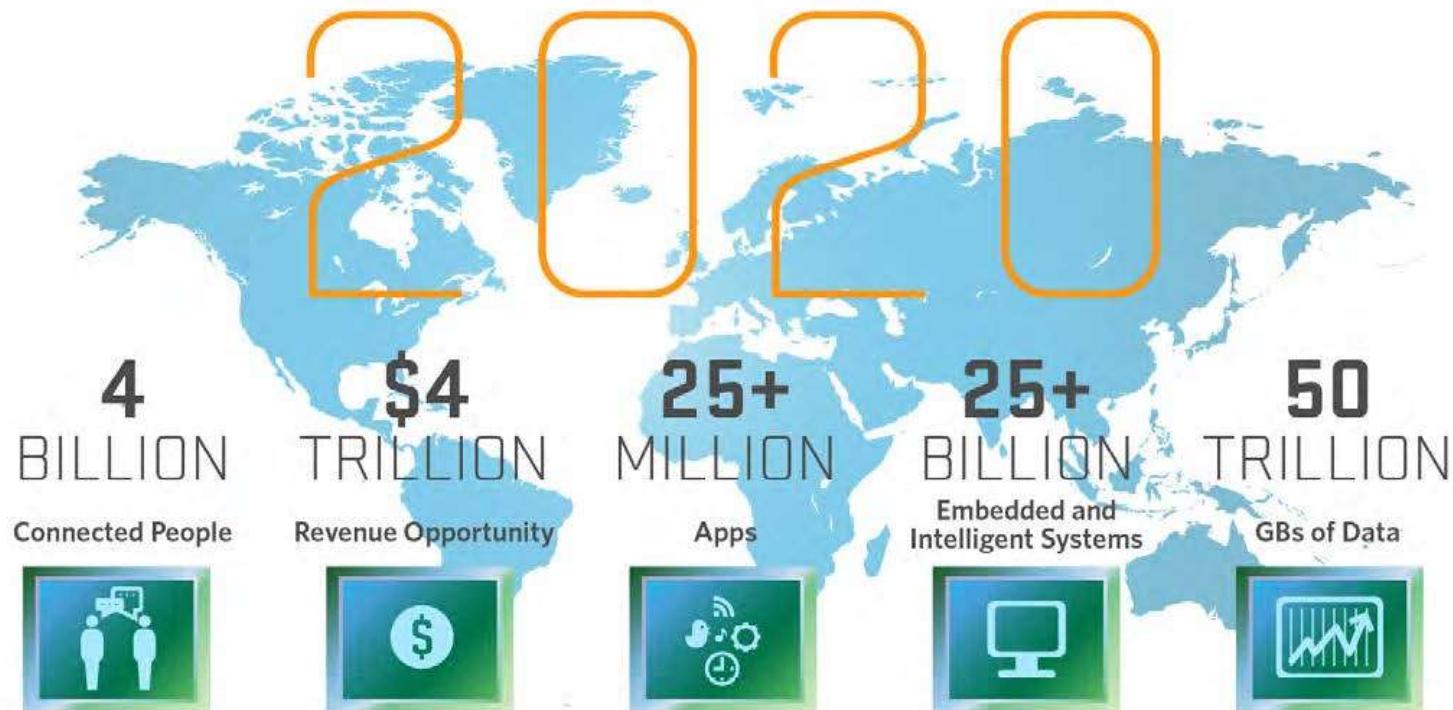
# Better Data, Better Government

- Quality and timely data are vital for enabling governments, international organizations, civil society, private sector and the general public to make informed decisions
- **Evidence based policy making**
- Quality of Statistics:
  - Accuracy
  - Relevance
  - Timeliness
  - Accessibility
  - Coherence
  - Interpretability

# Data Explosion

- Interactions of billions of people using computers, GPS devices, cell phones, and medical devices.
- online or mobile financial transactions, social media traffic, and GPS coordinates.
- “In the next five years, we’ll generate more data as humankind than we generated in the previous 5,000 years”. Eron Kelly, GM Microsoft

# Data Explosion



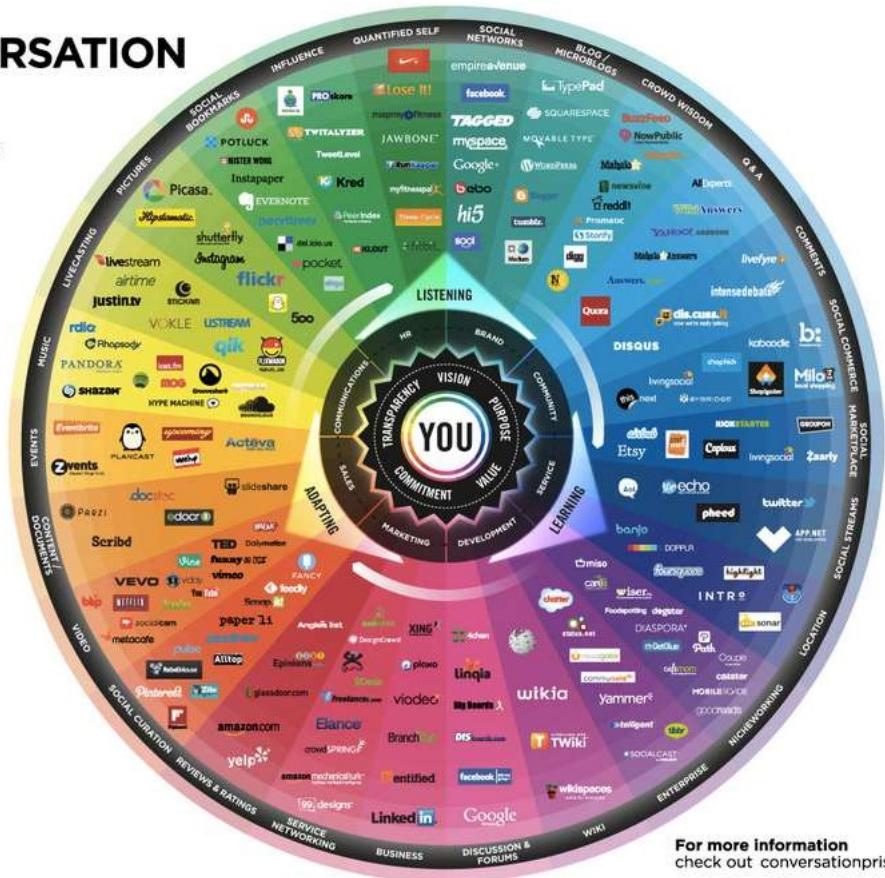
Source: Mario Morales, IDC



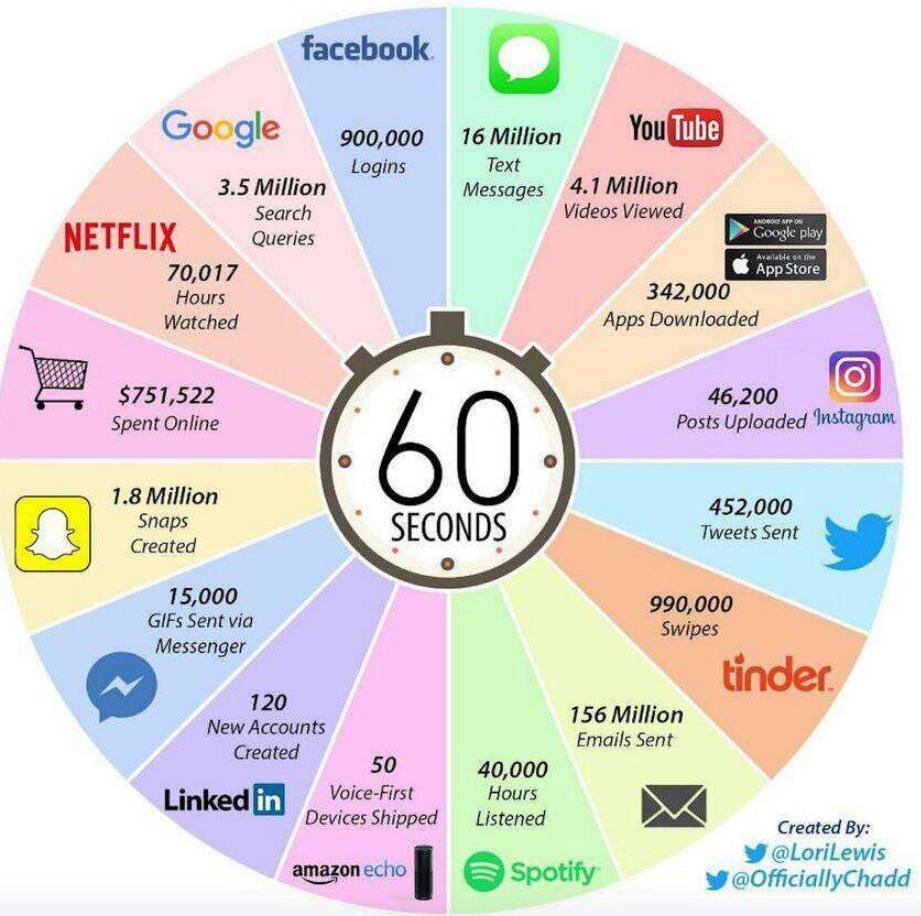
# What we “produce”?

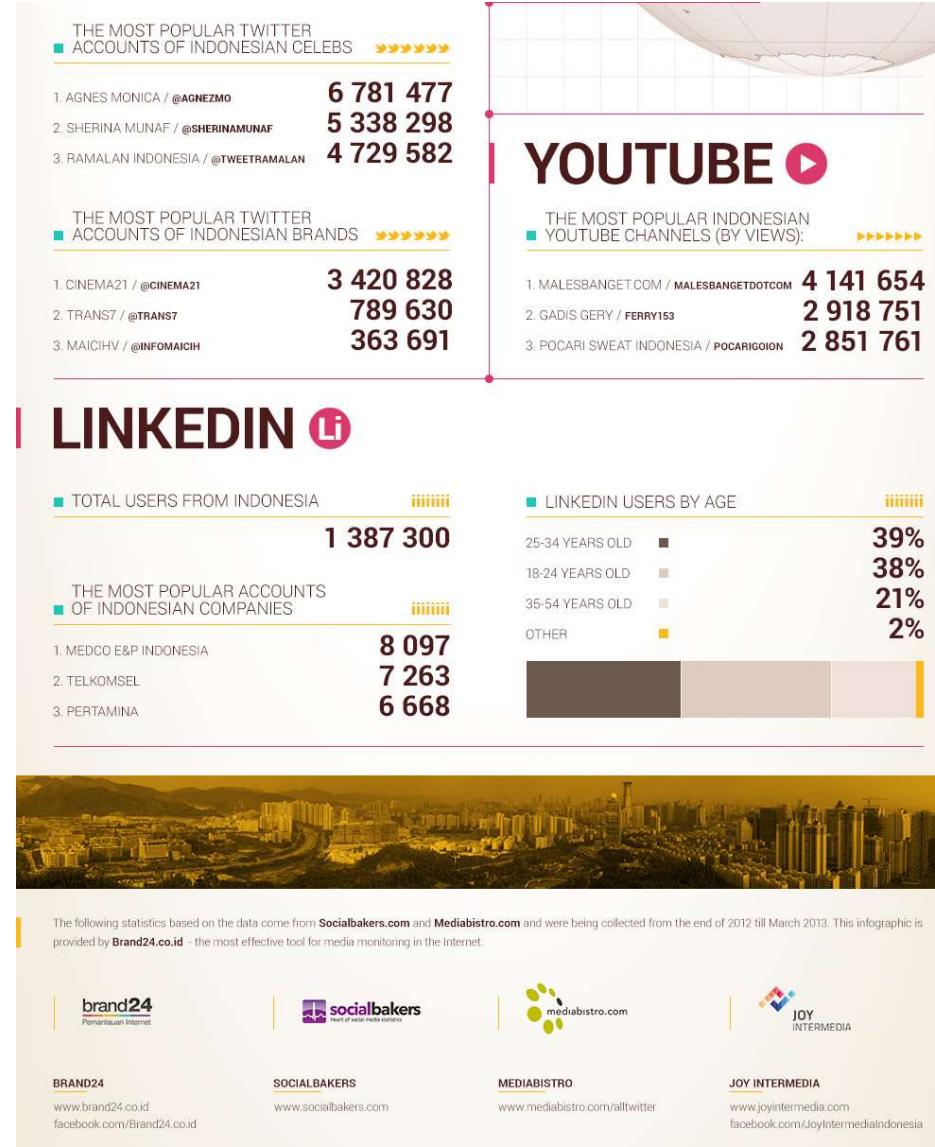
## THE CONVERSATION PRISM

Brought to you by  
Brian Solis & JESS3



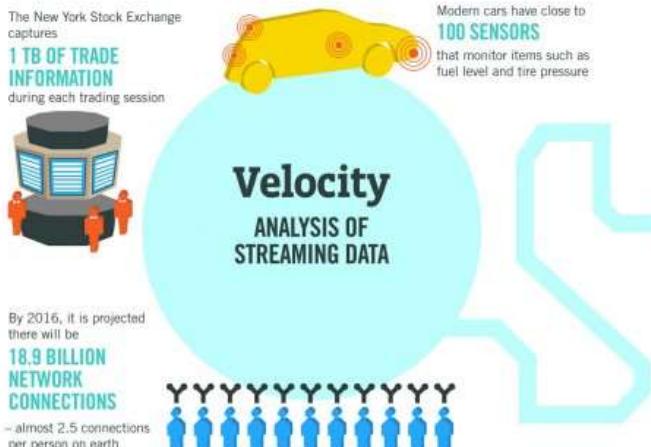
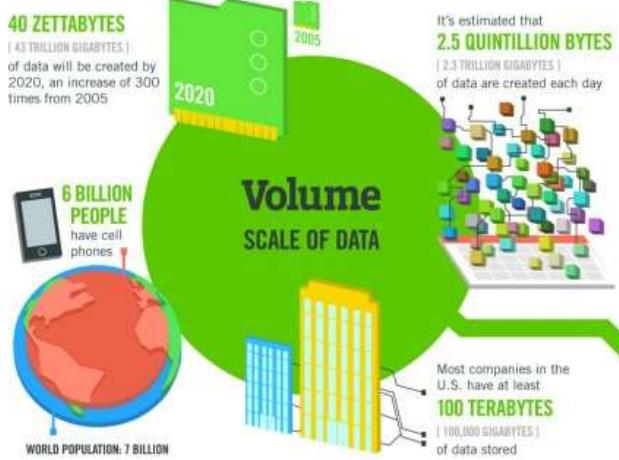
## 2017 This Is What Happens In An Internet Minute





# Internet of Things





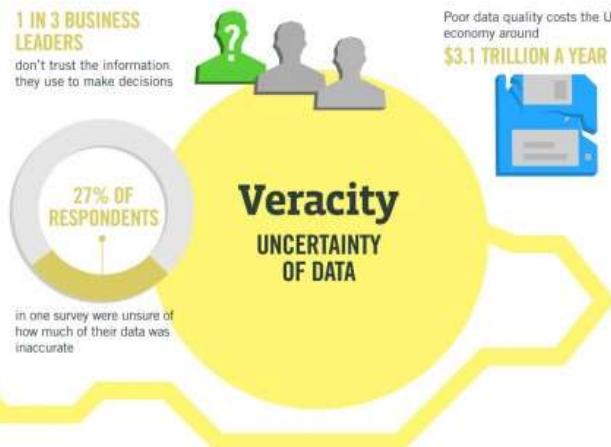
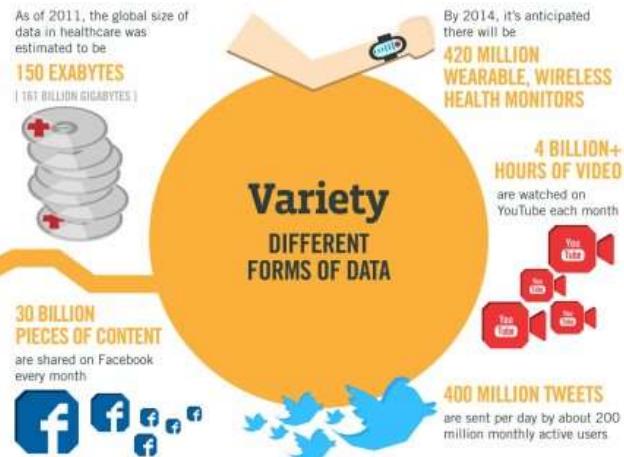
## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: Volume, Velocity, Variety and Veracity.

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015  
**4.4 MILLION IT JOBS**  
will be created globally to support big data, with 1.9 million in the United States



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MPTEC, Gartner

IBM

# Taxonomy BigData Sources

- **Exhaust data:** Passively collected data from people's use of digital services such as mobile phones, financial transactions or web searches.
- **Sensing data:** Actively collected data from sensors, e.g. in smart cities or from wearables and also through remote sensing and satellite images.
- **Digital content:** Open web content actively produced by people such as social media interactions, news articles, blogs or job postings. Unlike exhaust and sensing data is digital content intentionally edited by somebody, i.e. subjective or even deceptive, depending on the intentions of the author.

Letouzé (Data-Pop Alliance, 2015)

# Data Sources

<b>Exhaust data</b>	Mobile phone data
	Financial transactions
	Online search and access logs
	Citizen card
	Postal data
<b>Sensing data</b>	Satellite and UAV imagery
	Sensors in cities, transport and homes
	Sensors in nature, agriculture and water
	Wearable technology
	Biometric data
	Internet of Things (IoT)
<b>Digital Content</b>	Social media data
	Web scraping
	Participatory sensing / crowdsourcing
	Health records
	Radio content

# Measurement Revolution

**What People Do**

**What People Say**

# Mobile Positioning Data

- Location of Mobile Devices
- Statistical indicators can be generated:
  - The number of residences geographically distributed according to available accuracy;
  - The number of workplace, school, secondary home, and other regular locations;
  - Internal migration based on the change of the residences within the country;
  - Change of workplace over time;
  - Cross-border migration based on the regular travels between different countries;
  - Population grid statistics (1 km<sup>2</sup>);
  - Temporary population statistics
  - Assessing temporary population (hourly, daily, weekly, monthly, etc.);

# Twitter

- Trending Topics

**WHAT TO DO IN TWITTER?**

Age Group	Activities
20-25 YEARS OLD	- Re-tweet - Tweeting - Checking on trending topic - Checking on entertainment information account
26-30 YEARS OLD	- Tweeting - Re-tweet - Checking on trending topic - Checking on online news portal account
36-39 YEARS OLD	- Tweeting - Checking on trending topic - Checking on online news portal
40-45 YEARS OLD	- Favourite - Checking on trending topic - Checking on online news portal - Re-tweet - Voting to contest/TV program - Replying public figure/celebrity Twitter status - Tweeting

 | JAKPAT - Mobile Survey Platform Indonesia

# Crowdsourcing

- The process of getting work, funding or information, usually online, from a crowd of people.
- The word **Crowdsourcing** is a combination of Crowd & Outsourcing



# Official Statistics vs. Big Data

Official Statistics	Big Data
1. Structured and planned product	1. Largely unstructured unfiltered "data exhaust", i.e., by-product of digital products (transactions, web, social media, sensors)
2. Methodological and clear concepts	2. Poor analytics
3. Regulated	3. Unregulated
4. Macro-level but typically based on high volume primary data	4. Micro-level huge volume with high velocity (or frequency) and variety
5. High cost	5. Generally little, or no cost
6. Centralized; point in time	6. Distributed; real-time

Dr. Jose Ramon G. Albert, NSCB, Philippines

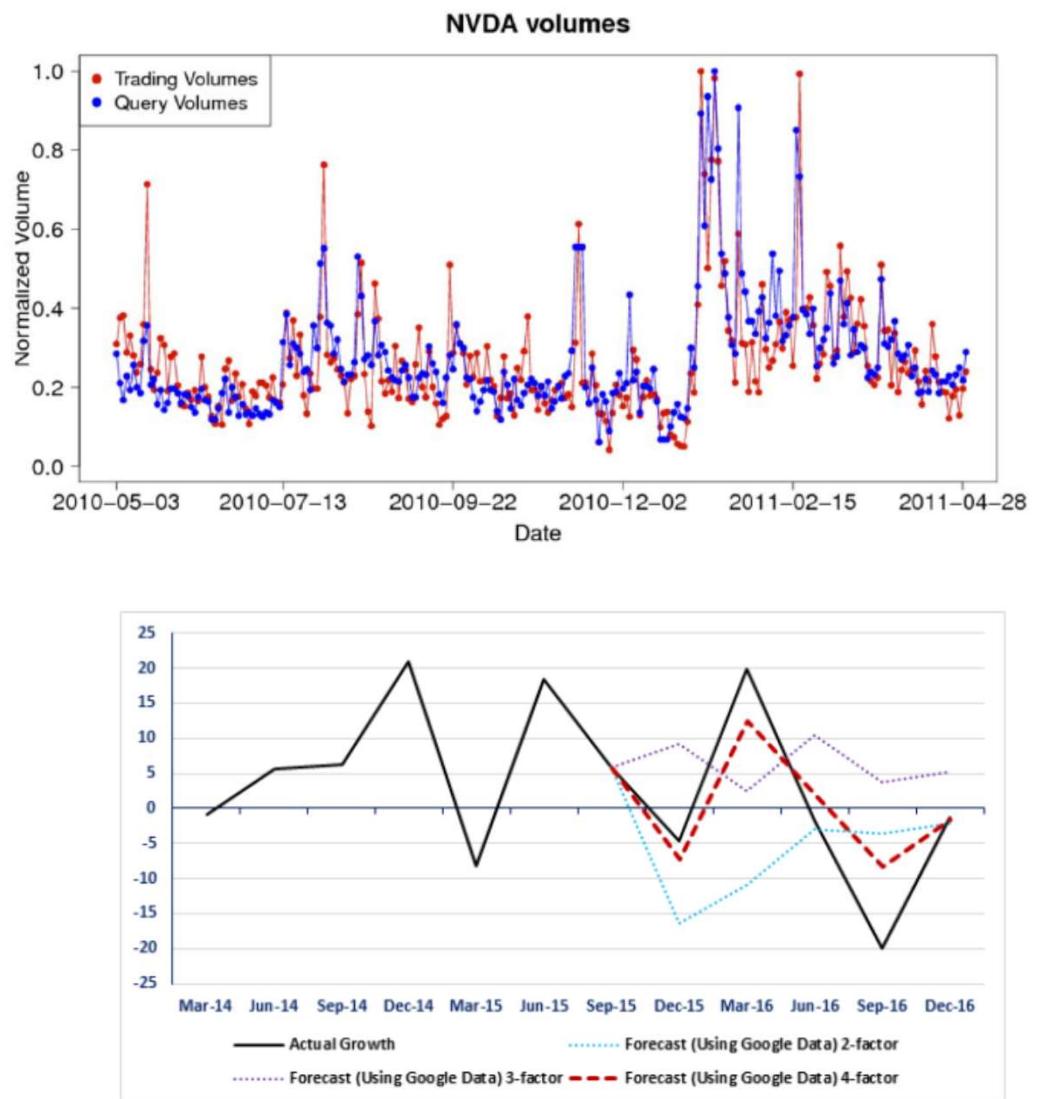
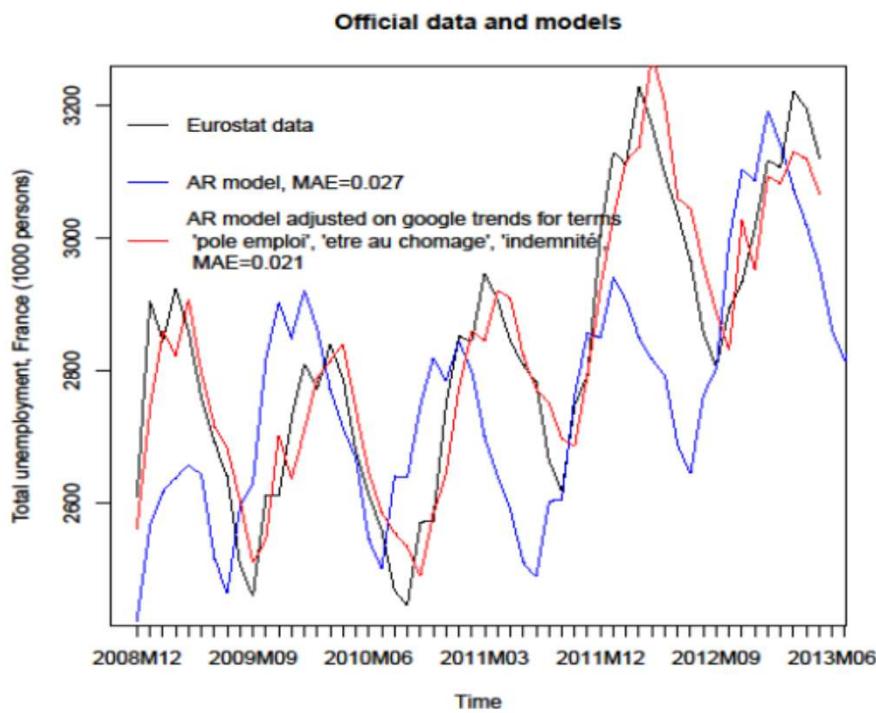
# Big Data Roles for Official Statistics

- Provide variables to help NSOs stratify better for sample surveys
- Improve sample survey estimates
- Help to compensate for nonresponse
- Help to check NSOs estimates
- Help to improve the frequency and timeliness of data releases
- Help to improve and provide more small-area estimates

*Cavan Capps and Tommy Wright, U.S. Census Bureau*

# Big Data for Public Policy: Examples

# Google Trend



ECONOMICS

## Predicting poverty and wealth from mobile phone metadata

Joshua Blumenstock,<sup>1,\*</sup> Gabriel Cadamuro,<sup>2</sup> Robert On<sup>3</sup>

Accurate and  
and economic  
enabling new  
sources of big  
be used to inf  
predicted attr  
distribution of  
composed of j  
and household  
and timely inf

## Mapping poverty using mobile phone and satellite data

## Mobile phone data could replace census questionnaires in the UK

The ONS used locations and timestamps to identify where people work and live.

tor A. Alegana<sup>1</sup>,

M. Iqbal<sup>6</sup>,

drew J. Tatem<sup>1,2,10</sup>



Matt Brian, @m4tt  
11.07.17 in Mobile

RESEARCH ART

## Mobile Phone Call Data as a Regional Socio-Economic Proxy Indicator

Sanja Šćepanović<sup>1,\*</sup>, Igor Mishkovski<sup>2</sup>, Pan Hui<sup>3</sup>, Jukka K. Nurminen<sup>1</sup>, Antti Ylä-Jääski

<sup>1</sup> Department of Computer Science, Aalto University, Helsinki, Finland, <sup>2</sup> Faculty of Computer Science a...  
Engineering, University Ss. Cyril and Methodius, Skopje, Macedonia, <sup>3</sup> Department of Computer Science  
and Engineering, Hong Kong University of Science and Technology, Hong Kong, China

4  
Comments

226  
Shares



telcor group research, Oslo, Norway

<sup>4</sup>School of Information, University of California, Berkeley, CA, USA

<sup>5</sup>Data Science Institute, Imperial College London, London, UK

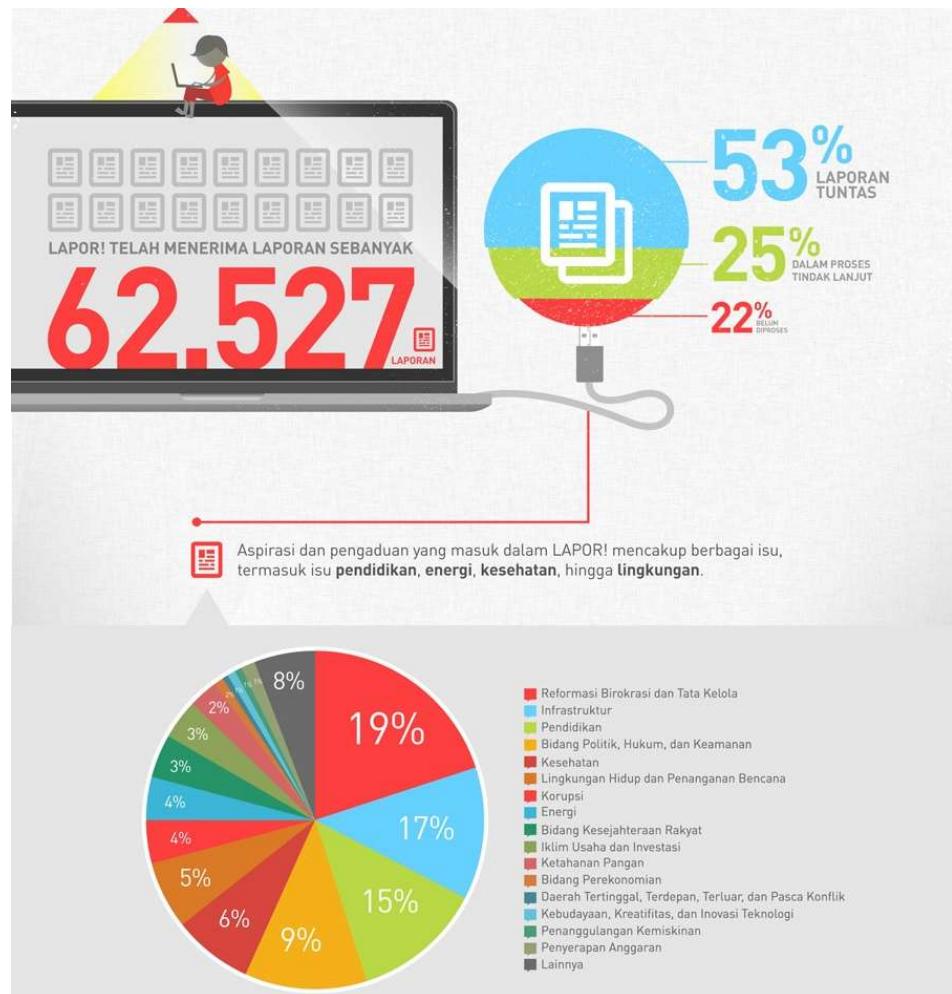
<sup>6</sup>Grameenphone Ltd, Dhaka, Bangladesh

<sup>7</sup>Public Health Sciences, Karolinska Institute, Stockholm, Sweden

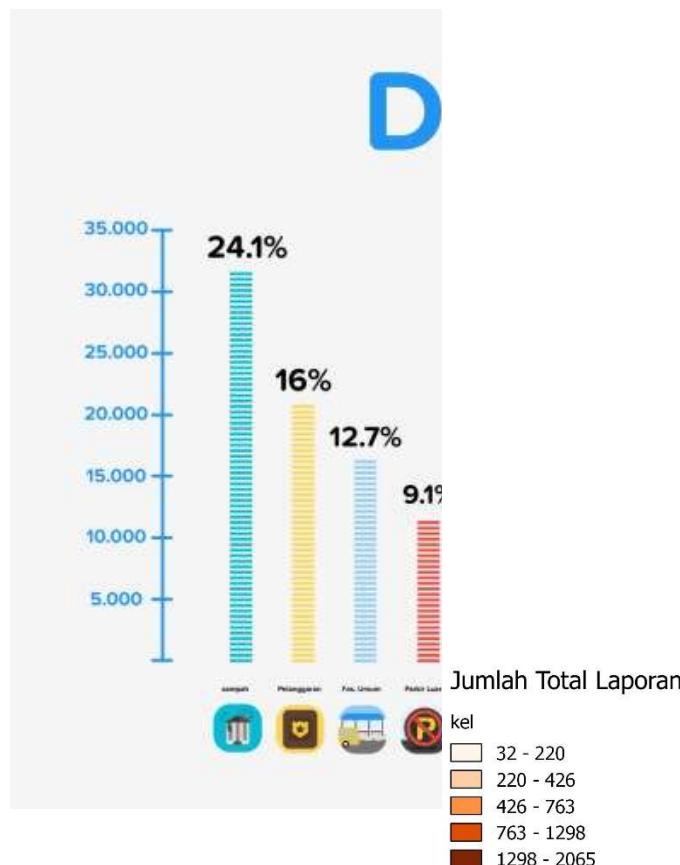
<sup>8</sup>College of Information System and Management, National University of Defense Technology,

g 44, Southampton, UK

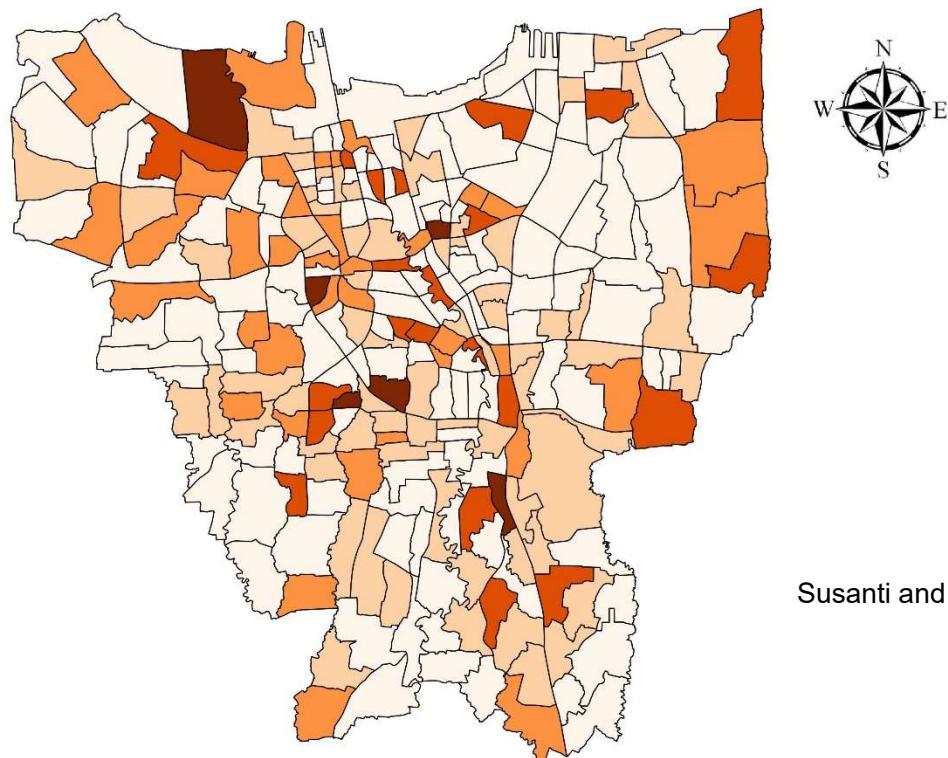
# Crowdsourcing



# Crowdsourcing



JUMLAH LAPORAN SAMPAH DI SETIAP KELURAHAN DI DKI JAKARTA TAHUN 2016

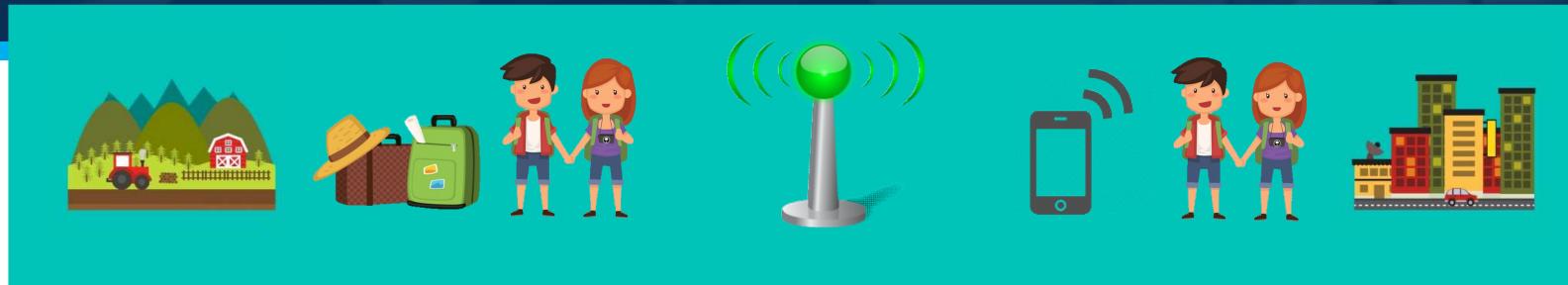


Susanti and Pramana 2017

# Web Crawling and Scraping

- Extract Information from Web
- Web Crawling is the process of locating information on World Wide Web(WWW), indexing all the words in a document, adding them to a database, then following all hyper links and indexes and adds that information also to the database
- Web scraping is the process of automatically requesting a web document and collecting information from it.

# Initiations of Big Data Use for Official Statistics by Politeknik Statistika STIS Jakarta



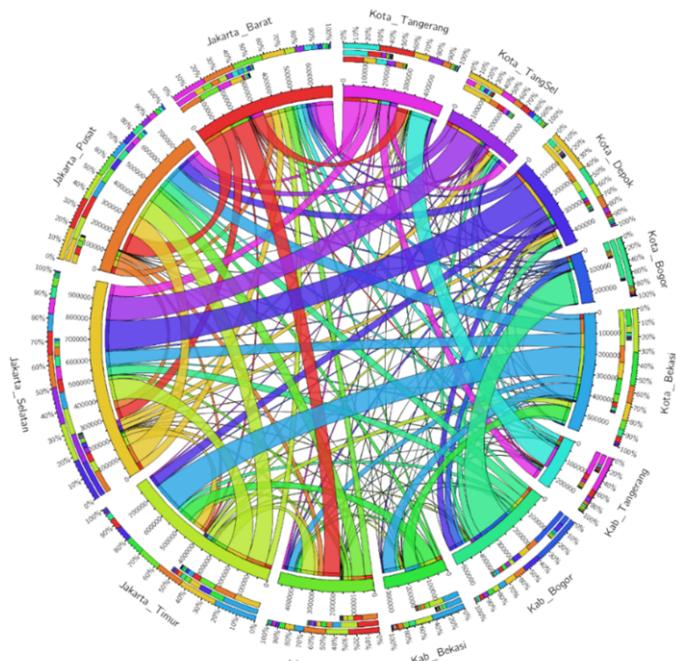
# Study Cases

- **Twitter Geotag for Predicting Commuting Patterns**
- **Price Nowcasting using crowdsourcing, online Shops and Instagram**
- **Happiness Index using twitter**
- **Mobility Behavior Tracking**
- **Job Vacancy Monitoring**
- **Tourism Promotion on Social Media and Online Media**

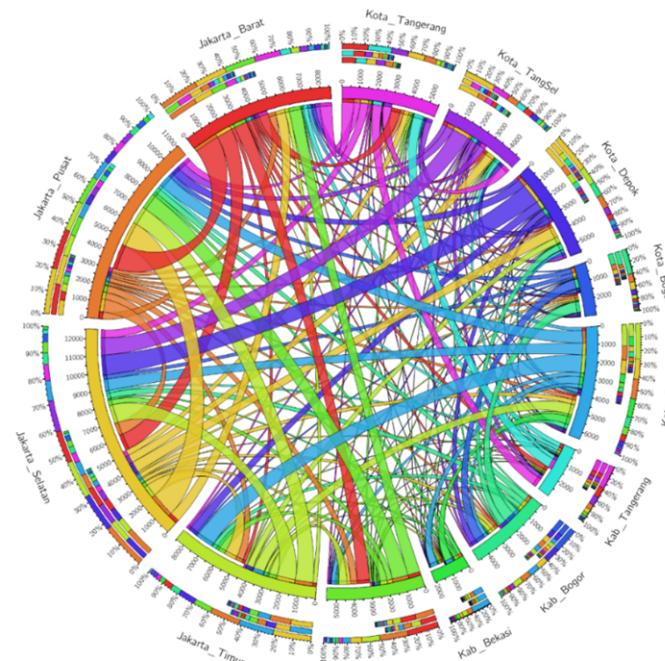
# Twitter Geotag for Predicting Commuting Patterns

- Collaboration with Global Pulse Jakarta
- Official Data: Jabodetabek Commuter Survey 2014
- Twitters (Februari 2014)
- Origin: Location of most tweets
- Destination: 2<sup>nd</sup> most tweets location

# Results



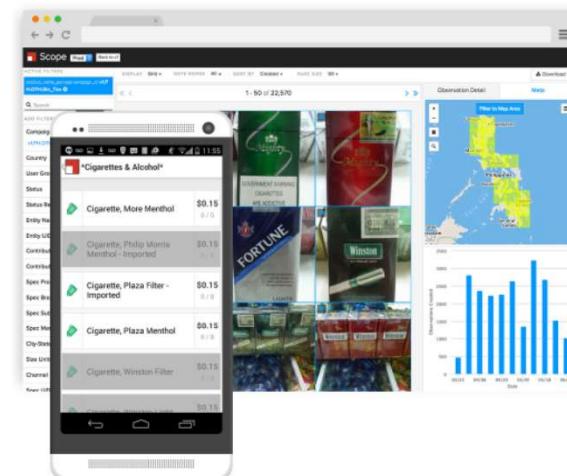
Commuter Survey

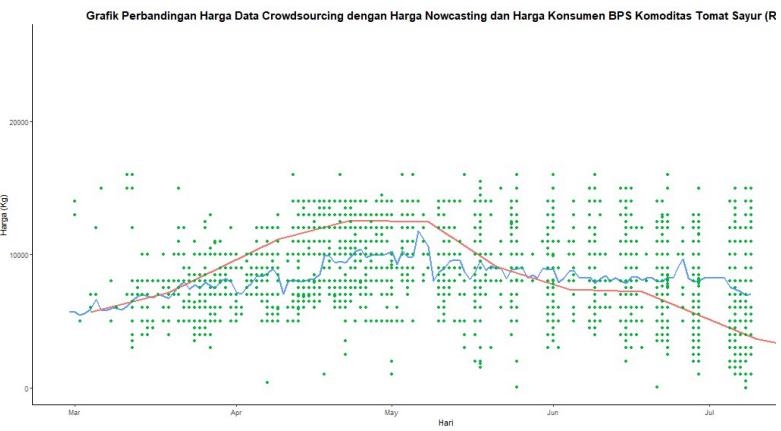
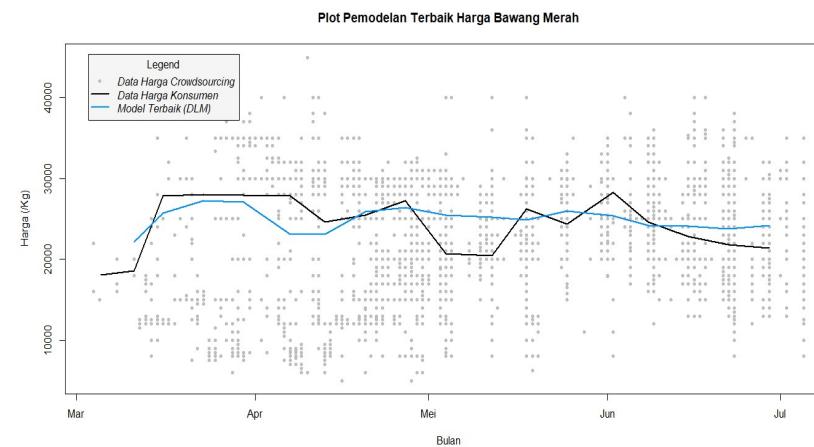
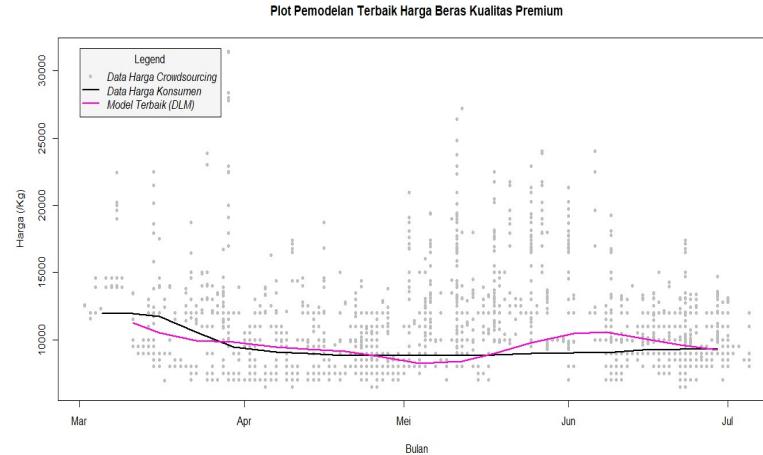
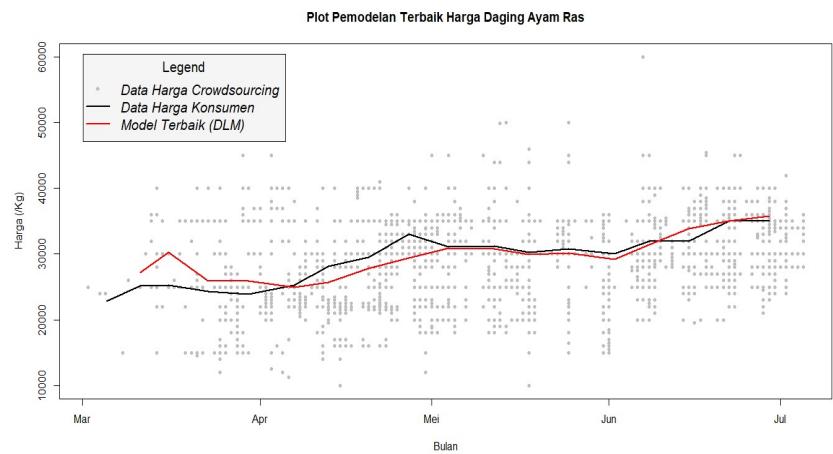


Twitter

## Crowdsource for Food Prices Nowcasting

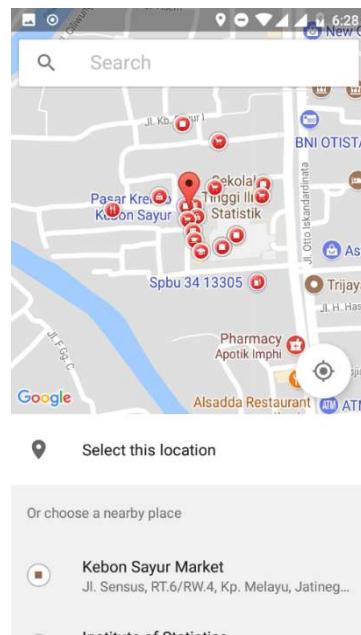
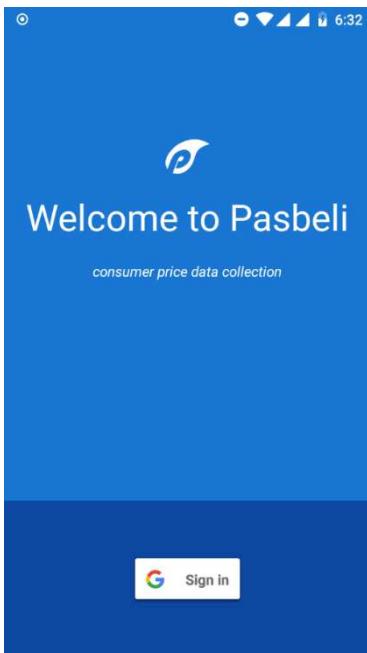
- Collaboration with Pulse Lab UN Jakarta
- Use crowdsourcing premise UN Food security project
- Locus: Kota Mataram, NTB
- Time: March– July 2015





# Crowdsourcing

## Pasbeli Apps: STIS Apps



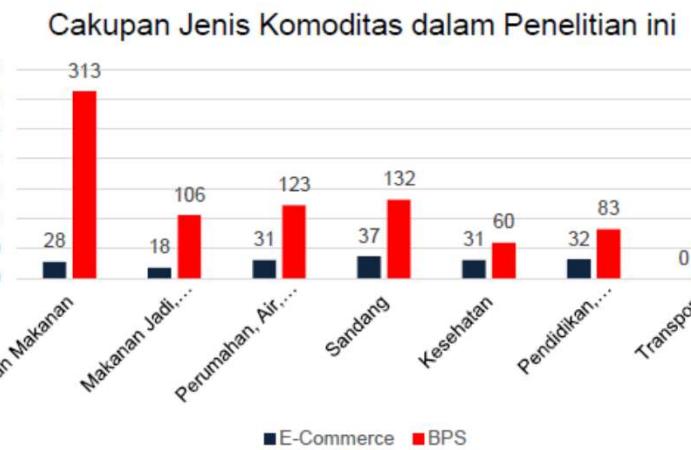
Setia Pramana

29



## Web Scraping: Online Shops

Online Shops	Total Commodities
Hypermart	52 products
KlikMart	75 products
Bhinneka	40 products
Elektronik City	17 products
Zalora	36 products
BerryBenka	25 products
Mothercare	2 products
Babyzania	4 products
Apotek Century	10 products
Pusat Kosmetik	5 products
Sephora	2 products
Stationary	6 products
Gramedia	3 products

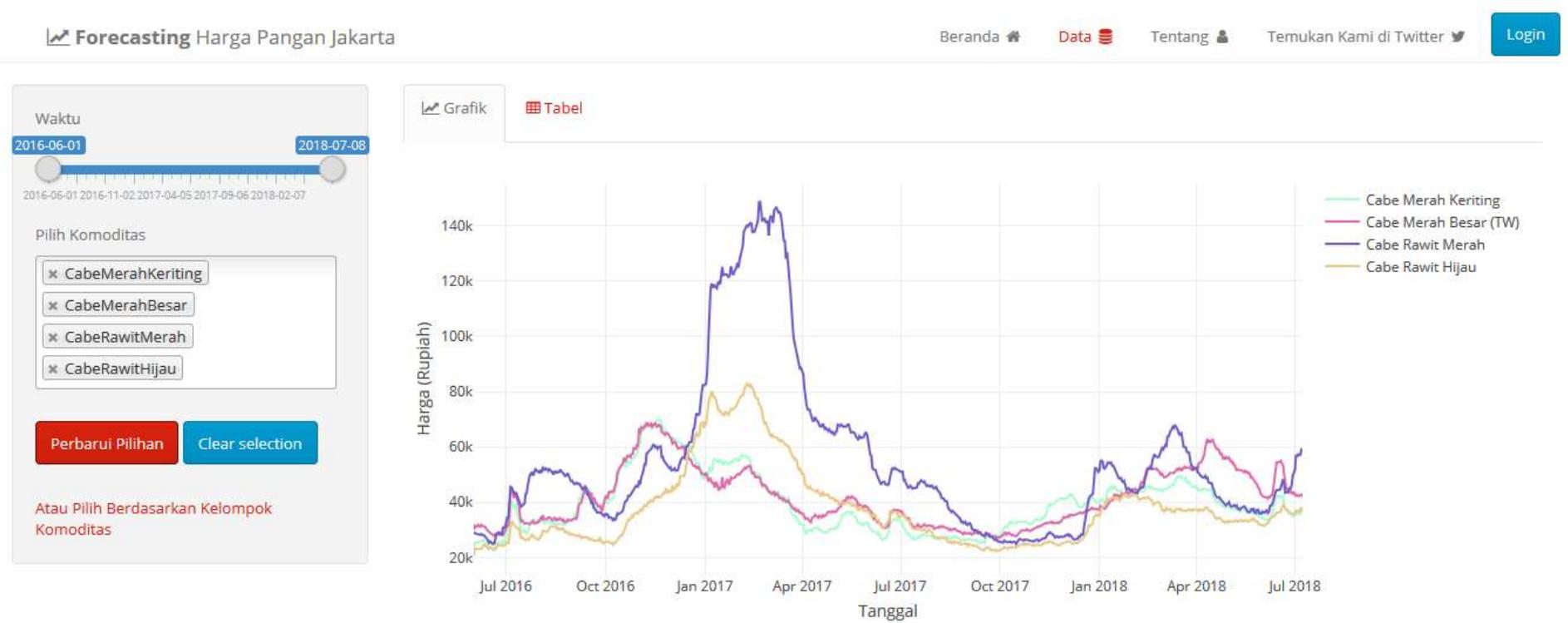




## Online Shop Commodities Prices

- Analysis is on progress
- Get the movement of consumer price
- Get the pattern of the changes of consumer price per commodity kind and per e-commerce
- Construct CPI by substituting the conventionally collected consumer price with e-commerce-based consumer price, then
- Comparing the survey-based CPI with e-commerce-based CPI

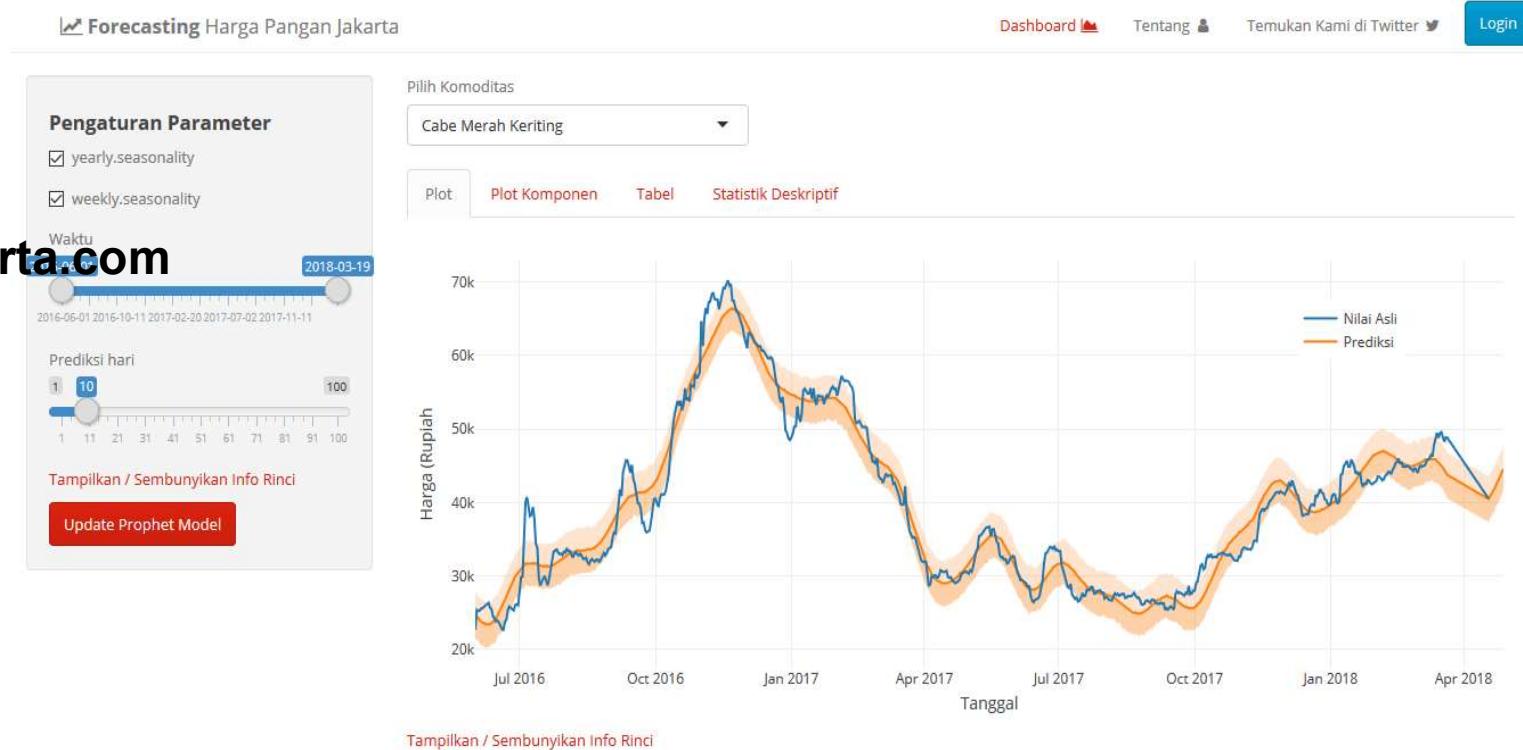
# Traditional Market Price Nowcasting



Muhammad Irsad Arief – 14.8260

# Traditional Market Price Nowcasting

Source of Scraping:  
<http://infopanganjakarta.com>





## Mobility Behavior Tracking

- Meili Travel Diary
- Collaboration with JUTPI phase 2 JICA
- KTH Stockholm
- Tracking people movement, speed, route...



## Mobility Behavior Tracking

- Record respondents travel behavior, including:
  - Trip routes
  - Destination (e.g.: Home, Office, Mall)
  - Trip activities (e.g.: Working, Shopping, Back to Home)
  - Transition Points (e.g.: Walking then using train)
  - First Level Transport modes (e.g.: Motorbike, Bus, Car)
  - Second Level Transport modes (e.g.: Commuter Line, Online Transport (Grab Bike/Car, Gojek))
  - Transport costs including: tariff, parking, and toll
- Forecasting respondents travel behavior (transport mode, activities) using machine learning algorithm



# Mobility Behavior Tracking

The figure consists of three screenshots from a mobile application for mobility behavior tracking.

**Screenshot 1: Map View**

A map of Jakarta showing a travel route. The route starts at "MULAI" (Start) and ends at "SELESAI" (End). The route is highlighted in red and shows a winding path through the city. The map includes various landmarks such as Mosjid Istiqlal, Taman Nasional, and Google Indonesia. A legend on the right indicates four map types: Road Map (selected), Satellite, Terrain, and Hybrid.

**Screenshot 2: Travel History Details**

This screen displays travel history details:

- Perjalanan sebelumnya berakhir pada Pukul 06:59 (Kamis, 02 Agustus 2018):**
  - Resume Lokasi Sebelumnya: 02 Agt. 08:59 - 02 Agt. 14:57, 7 jam 57 menit sebelum perjalanan ini.
  - Tempat: Pullman Jakarta Indonesia Thamrin Cbd
  - Aktivitas: Bekerja (Rapat)
- 14:57 (Kamis, 02 Agustus 2018) - Perjalanan Dimulai:**
  - Anda melakukan perjalanan sejauh 11 km dalam waktu 38 menit.
  - Dimulai pada: 02 Agt. 14:57
  - Berakhir pada: 02 Agt. 15:36
  - Menggunakan Moda: Mobil
  - Tipe: Taksi Konvensional
  - Tarif: 50000
- 15:36 (Kamis, 02 Agustus 2018) - Perjalanan Berakhir:**

**Screenshot 3: Segmentation Overview**

This screen shows a map of South Jakarta and Depok with a red travel route. It includes a summary of the trip segments:

- Anda berpindah moda transportasi pada titik ini (transition point)
- Titik ini merupakan lokasi tujuan perjalanan Anda (stop point)

Buttons for "Batal" and "Simpan" are visible.

# Job Vacancy Monitoring

Kane dan Alli Dikaitkan dengan Madrid, Ini Kata Luis Figo  
republika - 2018-02-28T07:18:00Z

tottenham hotspur liga spanyol la liga real madrid luis figo dele alli harry kane

Harga Minyak Dunia Turun Kali Pertama dalam Lima Hari  
republika - 2018-02-28T07:26:00Z

penurunan harga minyak minyak mentah dunia harga minyak naik harga minyak dunia

Domba Ini Berkeliaran dengan Isi Perut Terburai  
republika - 2018-02-28T04:55:00Z

australia hari ini domba dimangsa penggembala ternak australia serangan anjing liar  
australia plus abc

< 1 ... 7 8 9 10 11 ... 1349 >



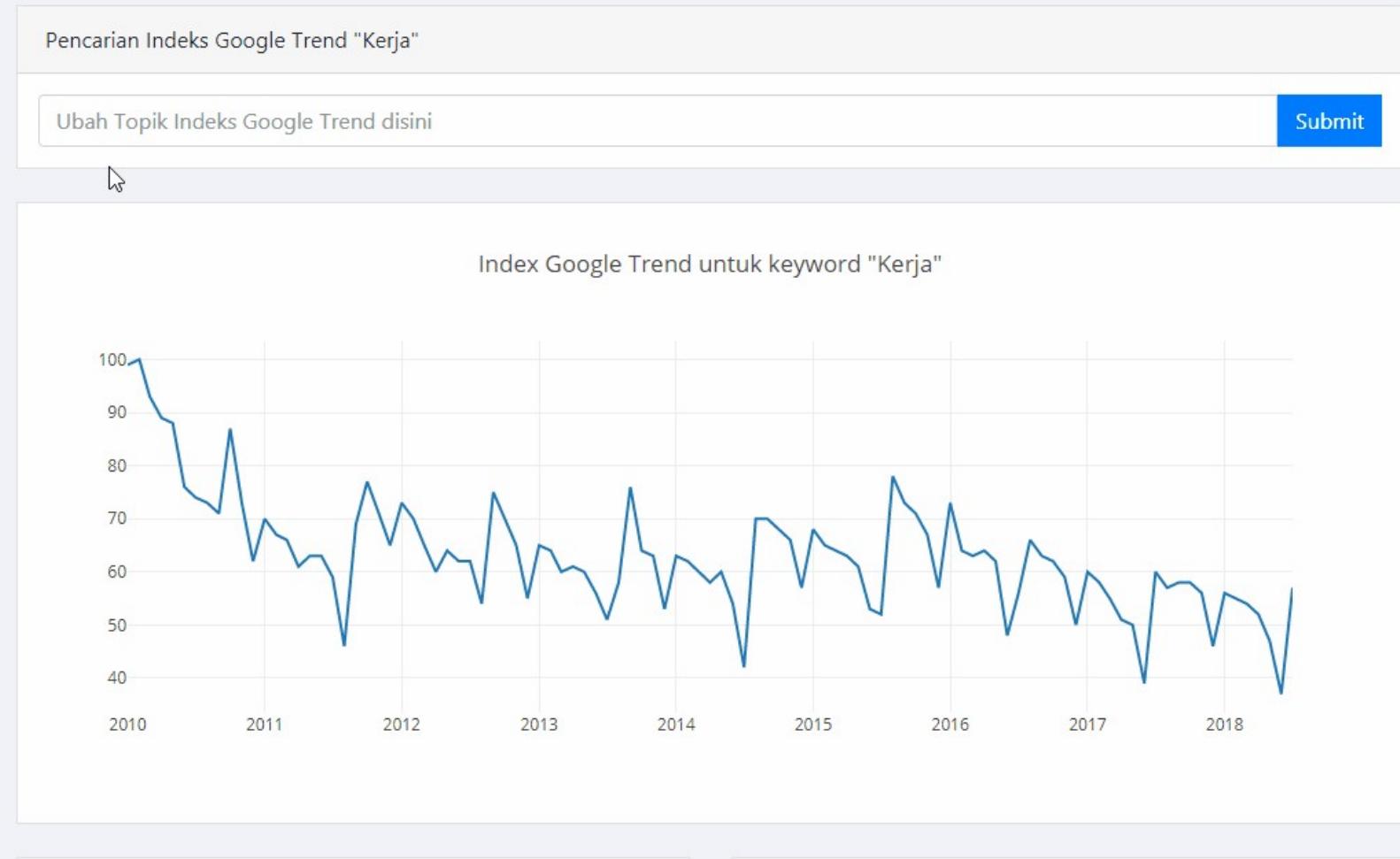
Persentase Kota untuk Lowongan yang Tersedia



TOP 5 Perusahaan Terfavorit  
(Jumlah Pelamar Terbanyak)



# Google trend



# Job Vacancy Monitoring: Twitter

Pencarian Twitter "Lowongan Kerja"

Ubah Topik Twitter disini

Submit

## Twitter dengan keyword "Lowongan Kerja"

Upssss ketuan lagi...kualitas kampret lsinya dikepalainya cuma Tai dan Tai...pantes mudah bicara dan menyebarkan h... <https://t.co/n8biTEBICF>

@rumahpositif99 Bekasi Timur, Indonesia 2018-08-14T17:11:17

Lowongan Kerja Supervisor Marketing Admin #lowker @ Jakarta, Indonesia  
<https://t.co/02gdTU7rVq>

@vendysjamsudin Setia Budi, Indonesia 2018-08-14T14:49:57

Lowongan RR Santa Tailor #Pekanbaru Agustus 2018: <https://t.co/A2aL6KaMsd>  
- Lowongan RR Santa Tailor #Pekanbaru Agu... <https://t.co/hT3m0e9Nh8>

@MermanGian Kerinci Kanan, Indonesia 2018-08-14T14:15:44

<https://t.co/WqBjJsBFRo> - Lowongan RR Santa Tailor #Pekanbaru Agustus 2018. RR Santa Tailor merupakan usaha konveks... <https://t.co/rAtRh8WaYP>

## Trending topic Indonesia

#ILCAntaraMaharDanPHP

Selamat Hari Pramuka

#AsianGames2018

#EmbraceTheEdgeOfPower

#SurpriseDeal

#JokowiRecoveryLombok

Farhat Abbas

Selasa 14 Agustus 2018

Kerajaan Ubur-Ubur

Turki

LIMA Basketball Nationals 2018

The Weeknd

#InvestasiUntukNegeri

#PertaminaPeduli

#KerjaNyataJokowi

#Recehkan17an

#DODIICEEDC07

# Scraping News

The screenshot shows the BITALISY web application interface. On the left, a dark sidebar contains the BITALISY logo and navigation links: Monitoring, Pencarian, Tabulasi, and Big Data. The main content area has a header "BITALISY | Scraping Berita". Below it is a form titled "Scraping Berita" with fields for selecting a website (dropdown menu showing "Detik") and choosing a date (text input field "mm/dd/yyyy"). A placeholder text "Silahkan Pilih Tanggal untuk Scraping Berita" is visible below the date input. A large blue button at the bottom right of the form says "Scraping Berita". To the right of the form is a "Log" section containing the message "Tidak ada scraping yang sedang berjalan". A cursor icon is positioned near the bottom center of the page.

**BITALISY** BITALISY | Scraping Berita

**Scraping Berita**

Pilih Website: Detik

Pilih Tanggal: mm/dd/yyyy

Silahkan Pilih Tanggal untuk Scraping Berita

Scraping Berita

**Log**

Tidak ada scraping yang sedang berjalan

# News Detail

**Bitalisy** 

BITALISY - BIG DATA ANALYSIS

- Monitoring
- Pencarian
- Tabulasi** >
- Big Data >

BITALISY | Tabulasi Berita | Detail

Penasaran dengan suatu topik? Cari disini  Submit

Tabulasi dari website "kompas" pada tanggal 01/02/2018

Show 10 entries

Search:

Tanggal ↑↓	Judul ↑↓	Berita ↑↓	Penulis ↑↓	Tag ↑↓
Feb. 1, 2018, 10:01 p.m.	Sandiaga: Kawasan Pasar Baru Ingin Dibuat seperti Boat Quay dan Clarke Quay di Singapura	JAKARTA, KOMPAS.com Wakil Gubernur DKI Jakarta Sandiaga Uno mengatakan pihaknya menerima usulan dari Armada Indonesia Kawasan Barat (Armabar) untuk membuat Sungai Ciliwung di kawasan Pasar Baru jadi taman kota percontohan. Wisata dan kegiatan komersil di Pasar Baru nantinya akan terintegrasi dengan sungai yang melintas di sampingnya. "Armada	Nibras Nada Nailufar	<a href="#">Sandiaga</a> <a href="#">Sungai Ciliwung</a> <a href="#">pasar baru</a>

# Jab vacancy

**Bitalisy** 

BITALISY | Tabulasi Lowongan Kerja

Penasaran dengan suatu topik? Cari disini  Submit

BITALISY - BIG DATA ANALYSIS

- Monitoring
- Pencarian
- Tabulasi >
- Big Data >

**Tabulasi Lowongan Kerja**

Show 10 entries 

[Copy](#) [CSV](#) [Excel](#) [PDF](#) [Print](#)

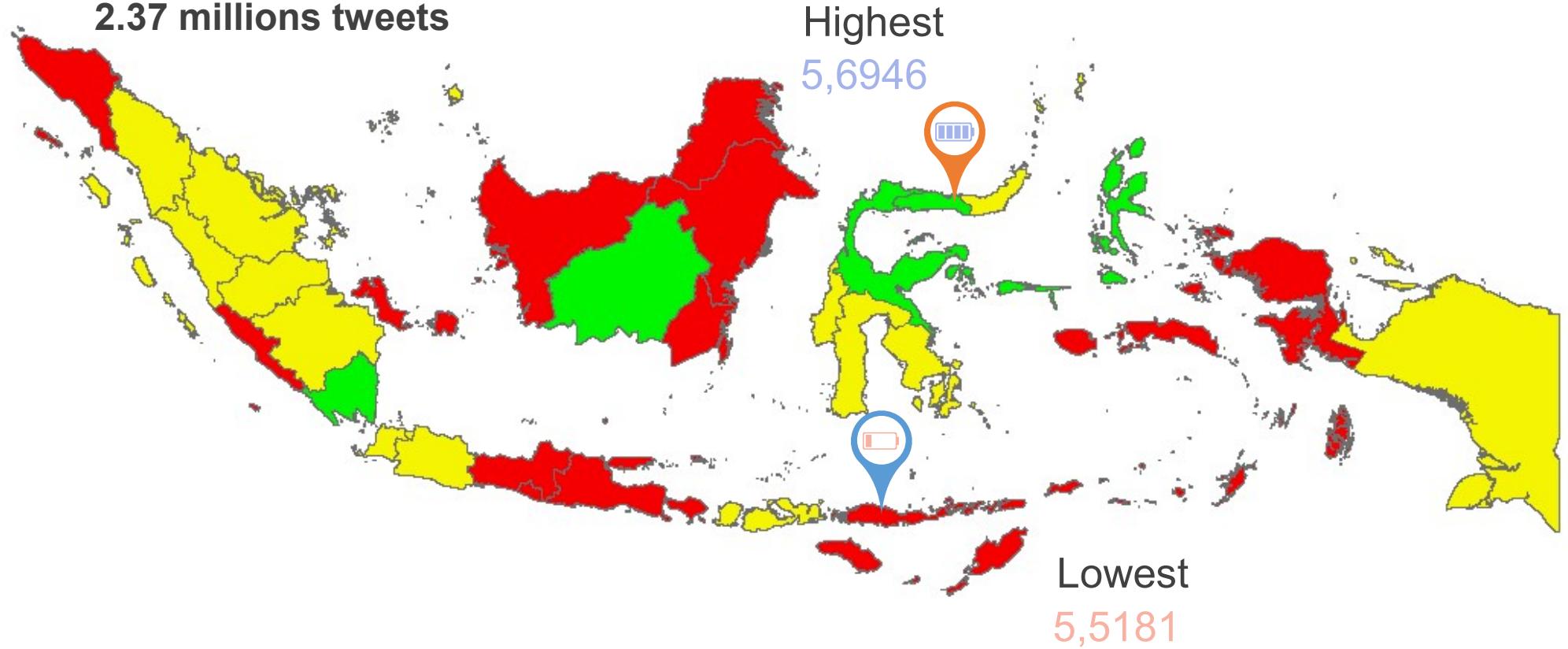
Search:

Diiklarkan sejak ↑↓	Ditutup pada ↑↓	Lowongan ↑↓	Perusahaan ↑↓	Gaji	Jumlah Pelamar ↑↓
Aug 01, 2018	Aug 31, 2018	TOUR Product Manager	Wisata Mega Utama PT	Gaji Dirahasiakan	0
Aug 01, 2018	Aug 31, 2018	Ticketing Staff	Multi MediacaPTa Mandiri Komunikasi PT	Gaji Dirahasiakan	128
Aug 01, 2018	Aug 16, 2018	Teknisi Restoran	Secret Recipe Indonesia PT	Gaji Dirahasiakan	0
Aug 01, 2018	Aug 31,	Staf Design Komunikasi	Beauty Kasatama CV	Fresh Graduate akan	398

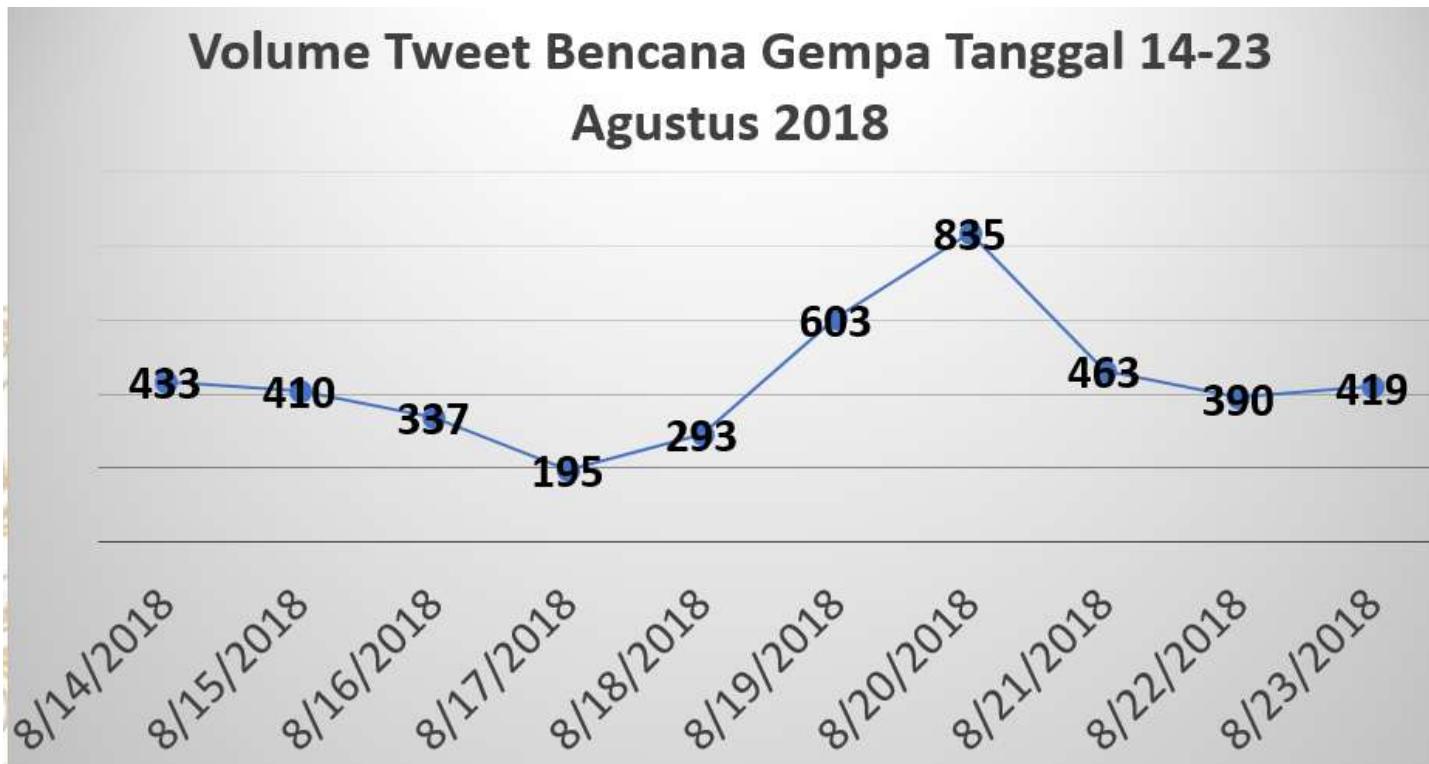
# Subjective Happiness Index: Twitter



Between April – July 2018,  
2.37 millions tweets



# Tweets: Disaster

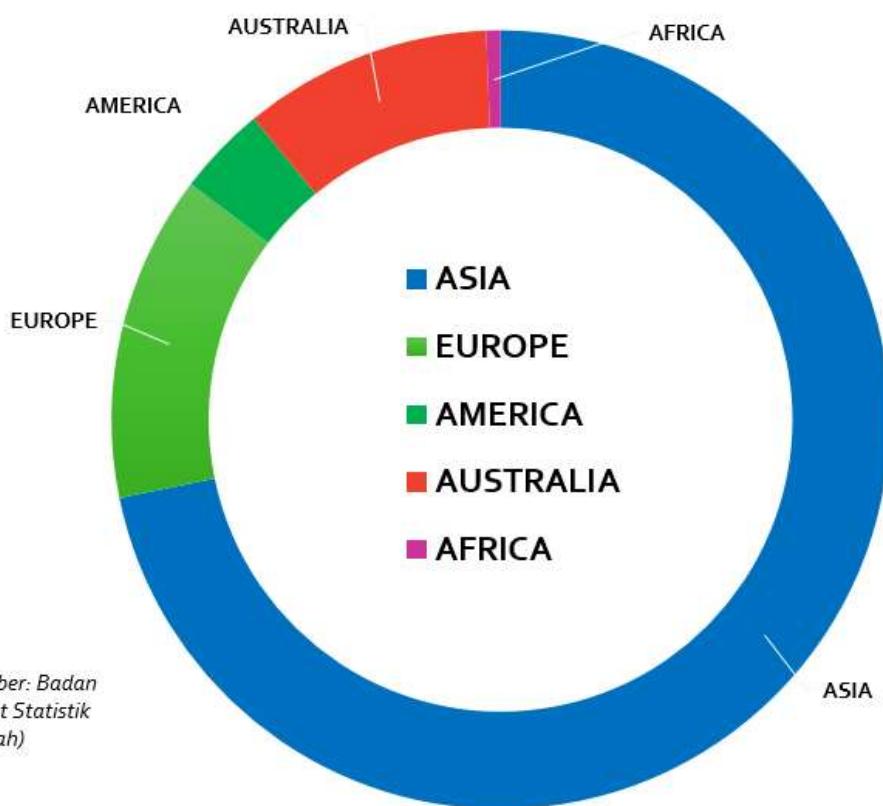


Lombok Earthquake 19 Agustus 2018.

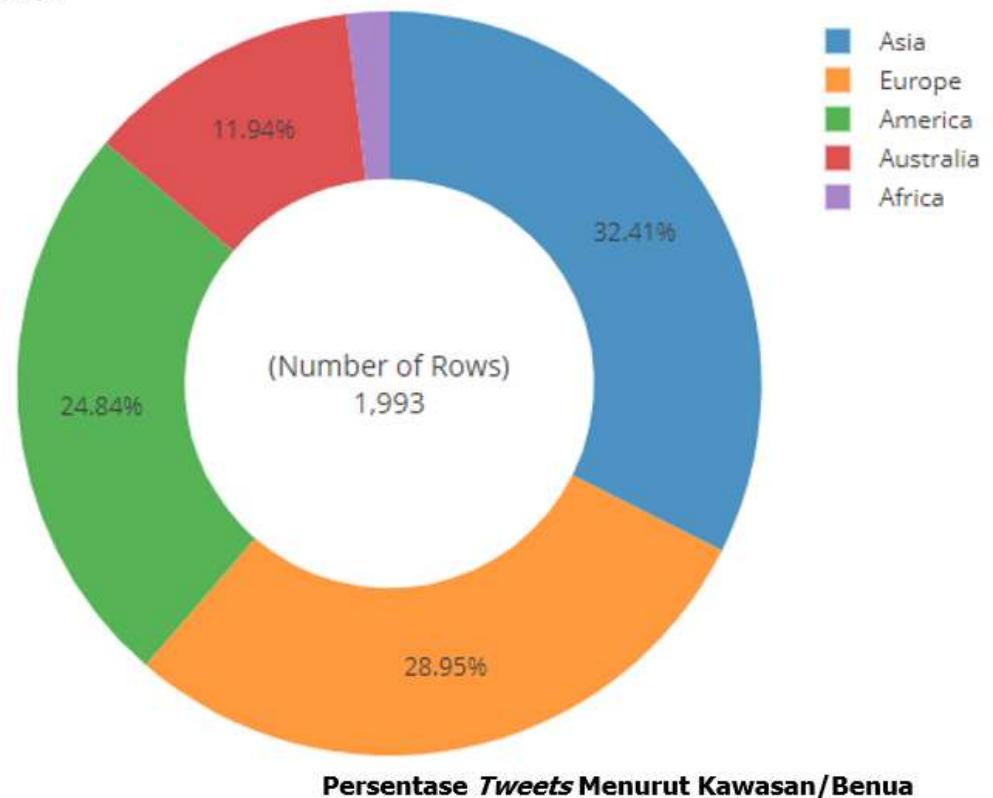


# Online News and Social Media Analysis of Indonesia Tourism Promotion

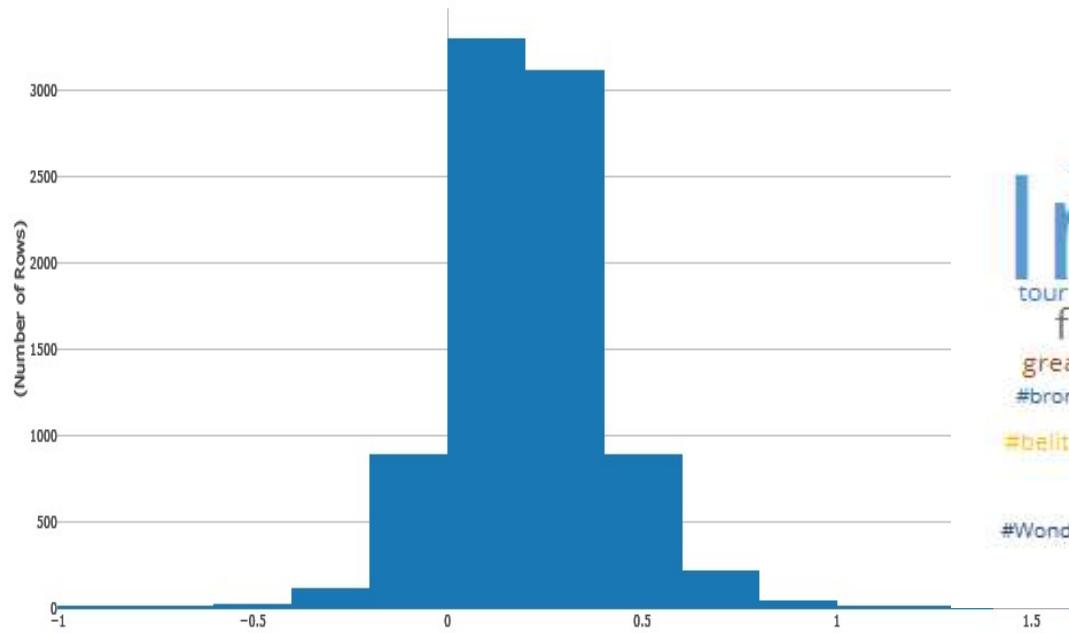
Percentase Wisatawan Manca Negara Menurut Kebangsaan Tahun 2017-2018



Sumber: Badan  
Pusat Statistik  
(diolah)



# Sentiment Analysis:

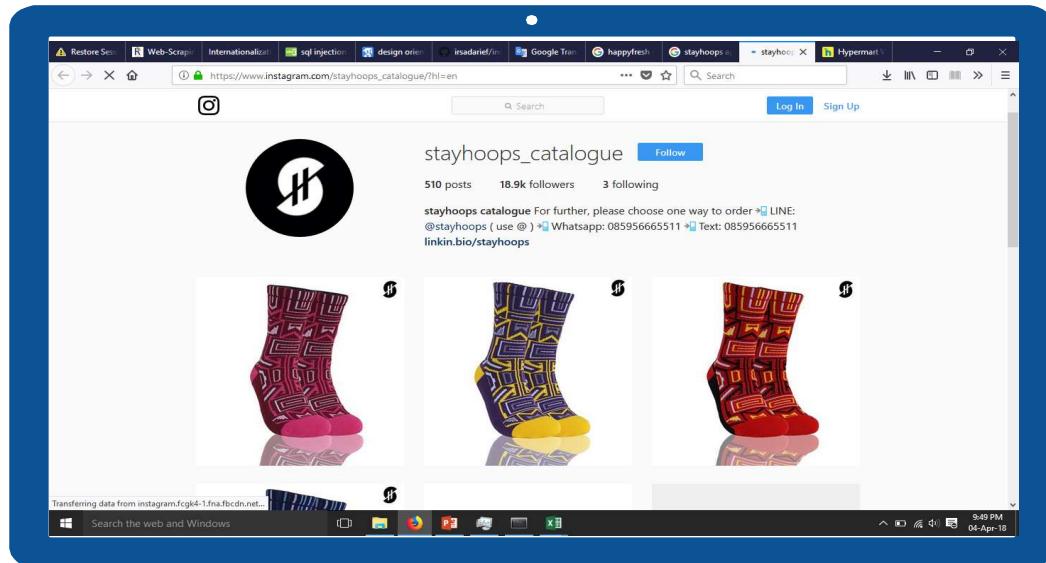


Most of Tweet are positives  
Further analysis are being performed



# Instagram Scrapping for Online Shopping

- 5 Million Instagram Users in Indonesia
- Get Price Products from Instagram



19k Followers  
currently provides +500 Products

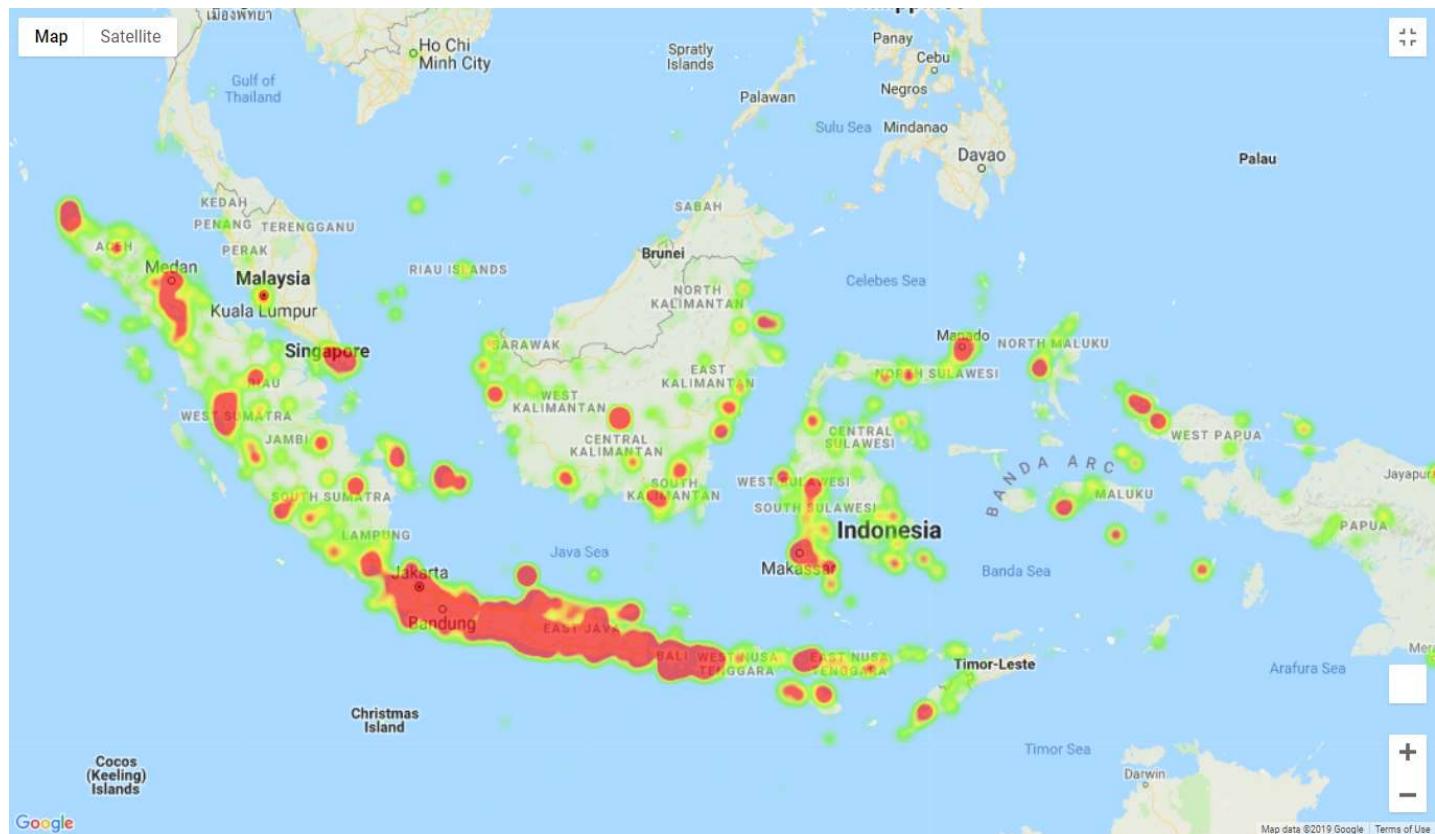
Original Text	Harga	Convert To
Ballerina short sleeve IDR.159.000 Available size: S, M, L, XL	159.000	159000
Grinding short sleeve IDR.159.000 Available size: S, M, L, XL	159.000	159000
Sleeveless Hoodie Blue Yellow IDR 275.000 Available size : S, M, L	275.000	275000

# Instagram Scrapping for Online Shopping: Challenges

1. Choose right accounts :
  - Some accounts display the price of the product in the picture, this will certainly make it difficult in taking the price data, therefore the selected account must write down the price of the product in its caption
  - Choose accounts that uses the same pattern in write down the product price's in each post
2. The products sold on instagram are mostly lifestyle products: shoes, clothes, hijabs, etc. whose prices tend to be stagnant and unchanged in long period of time
3. Each account has its own pattern in write down the product price (example: IDR 125K, Rp 125rb, Rp. 125.000, 125000, etc.)

# Instagram: Geotag Tourism

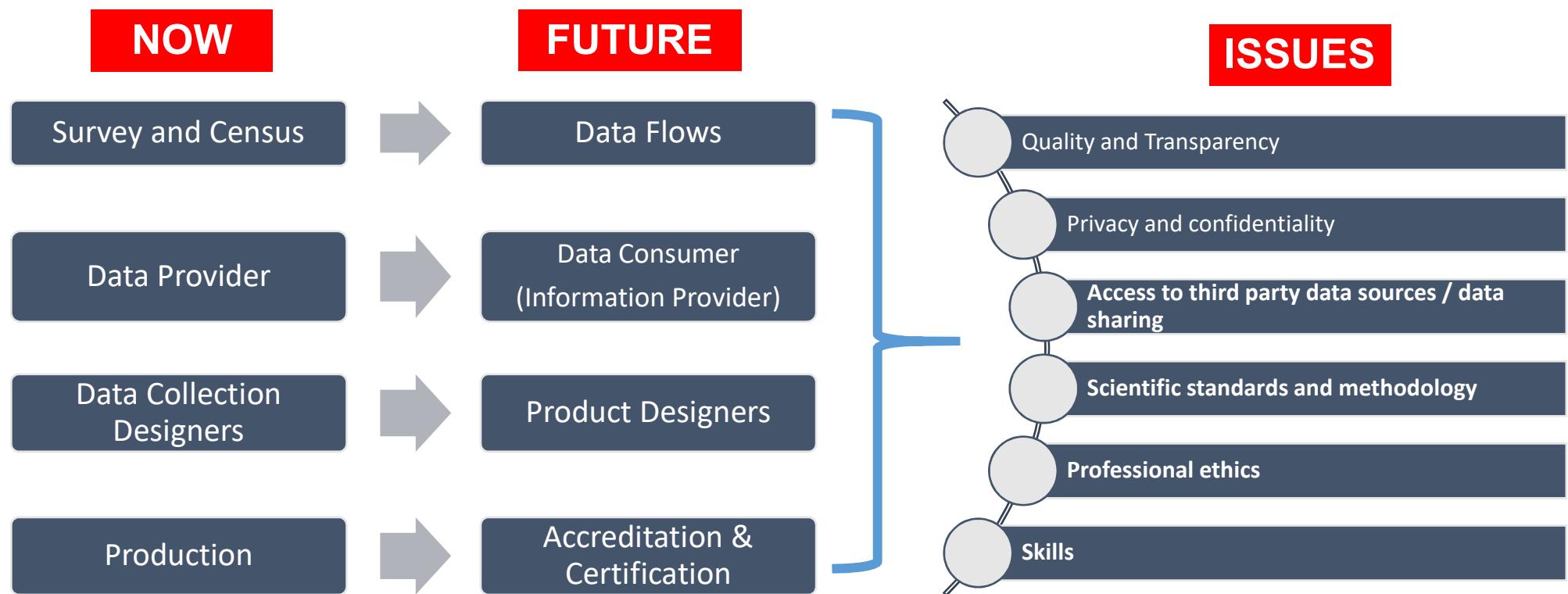
- #Wonderfulindonesia
- #Visitindonesia
- #Pesonaindonesia
- #Indonesiatourism
- #Exploreindonesia
- Total Images: 877.350
- With geo-tag : 225.328



# Current Projects

- Social Media Geospatial Data Analysis (Instagram) for:
  - Tourism Pattern in Indonesia
  - E-commerce (informal) analysis
- Automatic Web Scraping for:
  - Price nowcasting
  - Tourism Potential Analysis (Hotel Occupancy etc.)
- STIS data compilation and Big data analytics
- Study of the methodology of using Mobile Positioning Data (MPD) for commuter statistics in Indonesia.
- Image Analysis for Planting Pattern
- Study on big data use for:
  - Labor statistics
  - Unemployment
  - Telecommunication Statistics
  - Population census Combined method
  - Tourism Statistics

# Future National Statistics Office



# Challenges in Big Data Use for Official Statistics Production

- High dimensionality and extremely large size
- Possible coverage/selection bias
- Data accessibility, new legislation? Permission by public?
- Increase risk of data disclosure
- New Sampling Algorithm
- Heavy computation, new algorithm and analytic tools
- Integration of files from multiple sources in different format in different times
- Risk of data manipulation or sudden unavailability

# Summary

- Modernization of Statistical System is on progress...
- Still many approaches have to be explored and studied
- Contributions from **stakeholders** are needed