

Big Data: Introduction and Applications

August 20, 2015
HKU-HKJC ExCEL3 Seminar
Michael Chau, Associate Professor
School of Business, The University of Hong Kong

Big Data

- Ample opportunities for business organizations and governments to provide better services and gain managerial and strategic insights by gathering, cleaning, and analyzing these “Big Data”.

4

Over the past years...

- Data storage has grown exponentially
- Computation capacity has risen sharply
- Network bandwidth has increased greatly

2

What is Big Data?

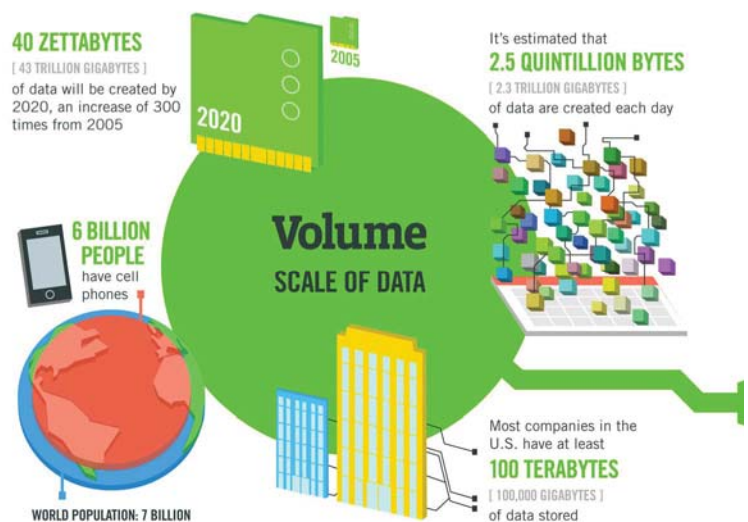
- It is not just big!
- Four Vs
 - Volume
 - Velocity
 - Variety
 - Veracity

5

Big Data

- Massive amount data are being generated at an unprecedented speed from various sources:
 - online transactions
 - mobile applications
 - sensors
 - images, audio, video
 - social media including blogs, weibos, facebook, and forums

3



The Data Size Is Getting Big, Bigger...

- Hadron Collider - 1 PB/sec
- Boeing jet - 20 TB/hr
- Facebook - 500 TB/day.
- YouTube – 1 TB/4 min.
- The proposed Square Kilometer Array telescope (the world's proposed biggest telescope) – 1 EB/day

Names for Big Data Sizes		
Name	Symbol	Value
Kilobyte	kB	10 ³
Megabyte	MB	10 ⁶
Gigabyte	GB	10 ⁹
Terabyte	TB	10 ¹²
Petabyte	PB	10 ¹⁵
Exabyte	EB	10 ¹⁸
Zettabyte	ZB	10 ²¹
Yottabyte	YB	10 ²⁴
Brontobyte*	BB	10 ²⁷
Gegobyte*	GeB	10 ³⁰

*Not an official SI (International System of Units) name/symbol, yet.

As of 2011, the global size of data in healthcare was estimated to be **150 EXABYTES** [161 BILLION GIGABYTES]

By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

4 BILLION+ HOURS OF VIDEO are watched on YouTube each month

30 BILLION PIECES OF CONTENT are shared on Facebook every month

400 MILLION TWEETS are sent per day by about 200 million monthly active users

Variety
DIFFERENT FORMS OF DATA

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session

Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure

Velocity
ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth

Google Analytics

1,912 people visited this site

Visits	Unique Visitors	Pages
21,281% (2,270 vs 3,211)	34.87% (3,023 vs 3,348)	5.75% (4,522 vs 4,273)
Pages / Visit	Avg. Visit Duration	Bounce Rate
26.97% (18,041 vs 1,027)	3.23% (00:01:53 vs 00:01:52)	18.52% (82,245 vs 91,205)

85.34% New Visitor
14.66% Returning Visitor

1 IN 3 BUSINESS LEADERS don't trust the information they use to make decisions

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**

27% OF RESPONDENTS in one survey were unsure of how much of their data was inaccurate

Veracity
UNCERTAINTY OF DATA

Mobile Devices

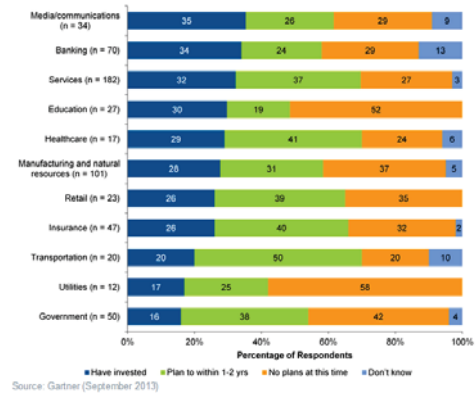
Social Media

- An important component of Big Data
 - Social networking sites
 - Blogs, microblogs
 - Online reviews, discussion forums
 - Online news aggregator
- People reveal themselves on social media
 - Demographics, preferences, habits, family ties, social ties
 - Images, videos
- Unstructured but rich in content
- Increasingly used in marketing research



13

Big Data Investment by Industry



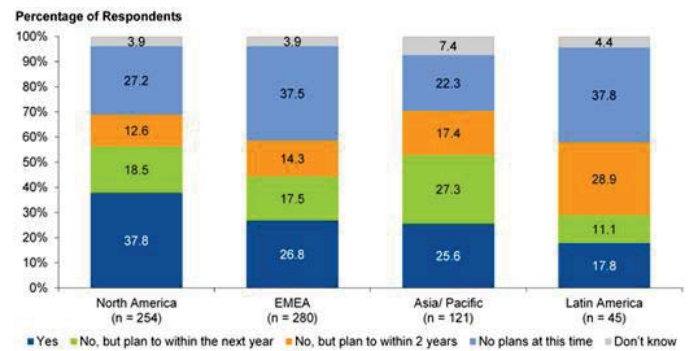
16

McKinsey estimates....



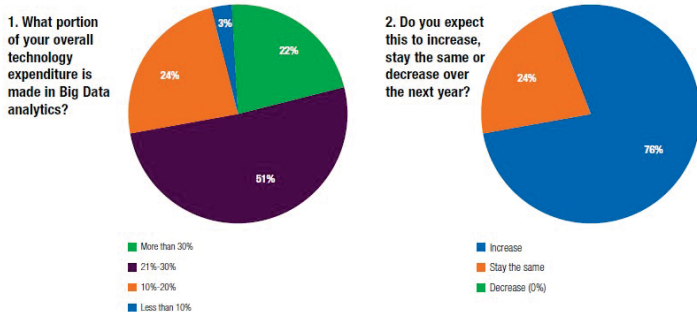
14

Big Data Investment by Region



17

Figure 1: Investments in Big Data analytics are strong



Source: Accenture "Industrial Internet Insights Report For 2015"

15

How governments see big data?

- USA: The White House Big Data R&D Initiative was launched in 2012
 - Extract knowledge and insights from large and complex collections of digital data
- UK: Formed the Alan Turing Institute for big data research
- South Korea: The Big Data Initiative was launched in 2011

18

How governments see big data?

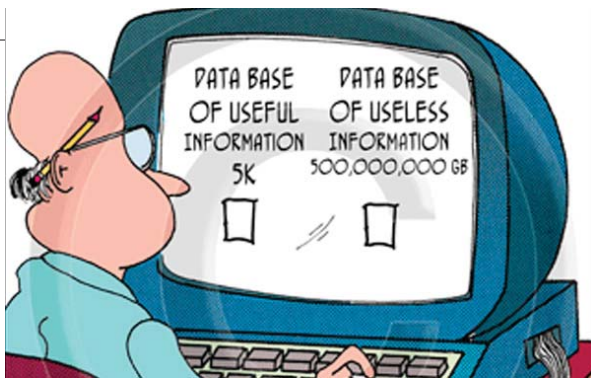
- United Nations' Global Pulse project
 - Analyze social media and big data for sustainable development and humanitarian action

19

Big Data Considerations

- You can't process the amount of data that you want to because of the limitations of your current platform.
- You can't include new/contemporary data sources (e.g., social media, RFID, Sensory, Web, GPS, textual data) because it does not comply with the data schema.
- You need to (or want to) integrate data as quickly as possible to be current on your analysis.
- You want to work with a schema-on-demand data storage paradigm because the variety of data types.
- The data is arriving so fast at your organization's doorstep that your analytics platform cannot handle it.
- ...

22



Source: The Storage Alchemist²⁰

Critical Success Factors for Big Data Analytics

- A clear business need (alignment with the vision and the strategy)
- Strong, committed sponsorship (executive champion)
- Alignment between the business and IT strategy
- A fact-based decision-making culture
- A strong data infrastructure
- The right analytics tools
- Right people with right skills

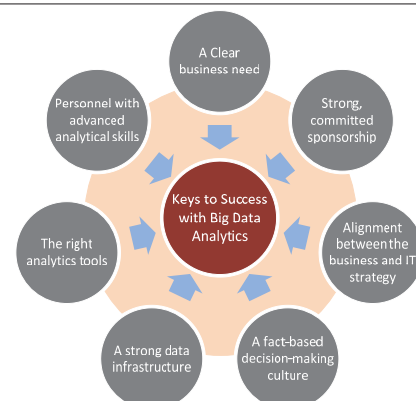
23

Fundamentals of Big Data Analytics

- Big Data by itself, regardless of the size, type, or speed, is worthless
- Big Data + "big" analytics = value (the 5th V)
- With the value proposition, Big Data also brought about big challenges
 - Effectively and efficiently capturing, storing, and analyzing Big Data
 - New breed of technologies needed (developed or purchased or hired or outsourced ...)

21

Critical Success Factors for Big Data Analytics



24

Models and Technologies

- Traditional
 - Data mining, data warehouse
- Big Data Analytics
 - Unstructured data from multiple sources
 - Real-time analytics
 - Complex statistical analysis
 - Distributed computing

25

Big Data Vendors

- Big Data vendor landscape is developing very rapidly
- A representative list would include
 - Cloudera - cloudera.com
 - MapR – mapr.com
 - Hortonworks - hortonworks.com
 - Also, IBM (Netezza, InfoSphere), Oracle (Exadata, Exalogic), Microsoft, Amazon, Google, ...

Software,
Hardware,
Service, ...

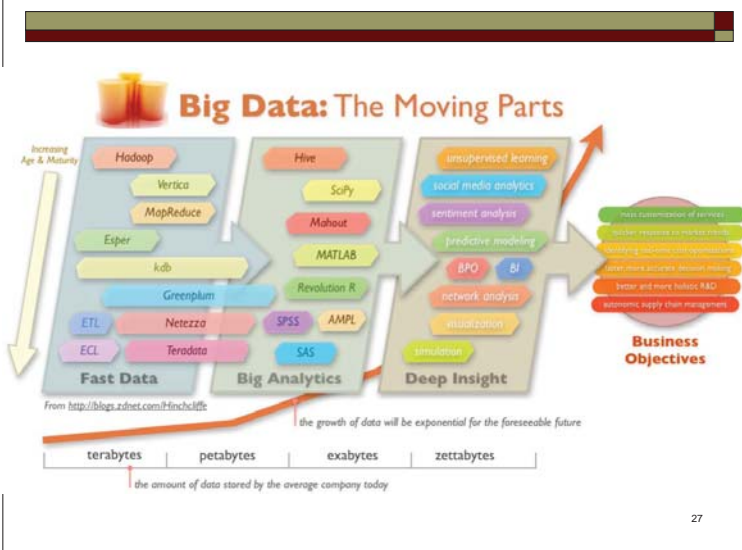
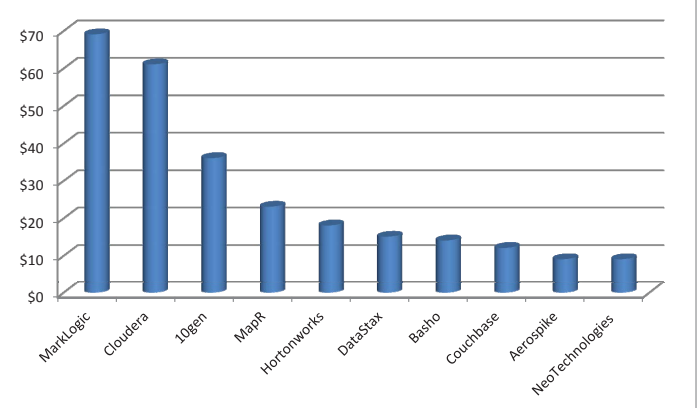
28

Models and Technologies

- Social media analytics
 - Text and sentiment analysis
 - Social network analysis
- Predictive modeling
 - Statistical method
 - Artificial intelligence/data mining
- Visualization

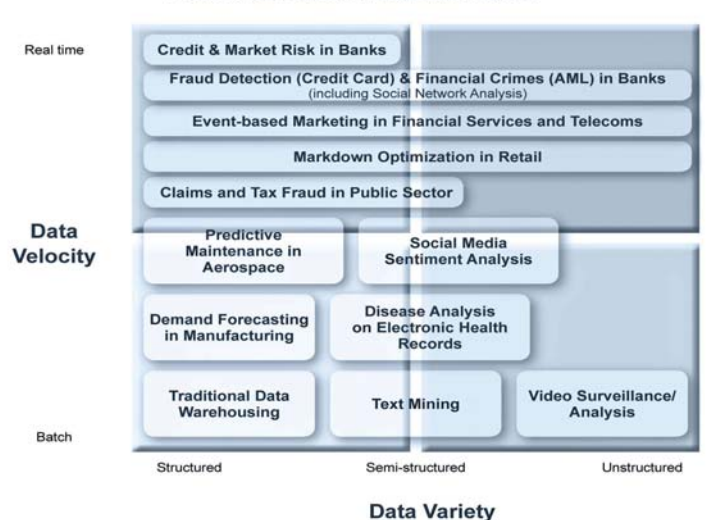
26

Top 10 Big Data Vendors with Primary Focus on Hadoop



27

Potential Use Cases for Big Data Analytics



Example: Amazon

- ❑ Predicts what customers will buy and manage their inventory
- ❑ “Customers who bought this item, also bought these...”
- ❑ Anticipatory Shipping: shipping an item to a customer in anticipation that this customer will order that product – will it work?

31

Example: Hewlett-Packard

- ❑ Analyzes data of 330,000 employees to predict who have a high risk of leaving the job
- ❑ Results in an estimated saving of \$300 million

34

Example: Google Flu Trend

- ❑ Google Flu Trend: Predicts influenza breakout based on the occurrences of relevant search terms in search data
- ❑ Successfully predicts regional outbreaks of flu up to 10 days before they were reported by the CDC (Centers for Disease Control and Prevention).

32

Example: Germany Soccer

- ❑ Match Insights: Collects and analyzes massive amounts of player performance data (including video data)



Example: Macy's

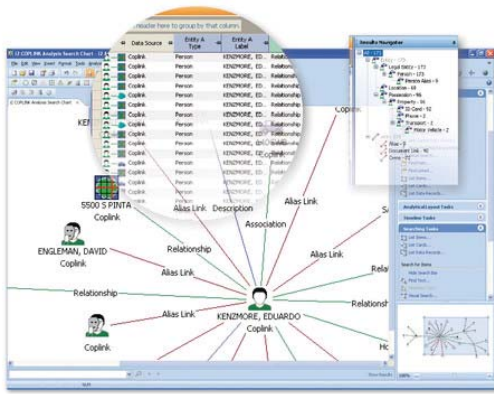
- ❑ Analyzes a vast amount of customer data ranging from visit frequencies and sales to style preferences and personal motivations.
- ❑ Adjusts pricing in near-real time for millions of items, based on demand and inventory
- ❑ Sends targeted and customized direct mailings to customers

33

Example: US Law Enforcement

- ❑ Challenges
 - Isolated databases
 - Lack of analytics capability
- ❑ COPLINK system
 - Implemented in many police departments in the US
 - Database linking
 - Association mining

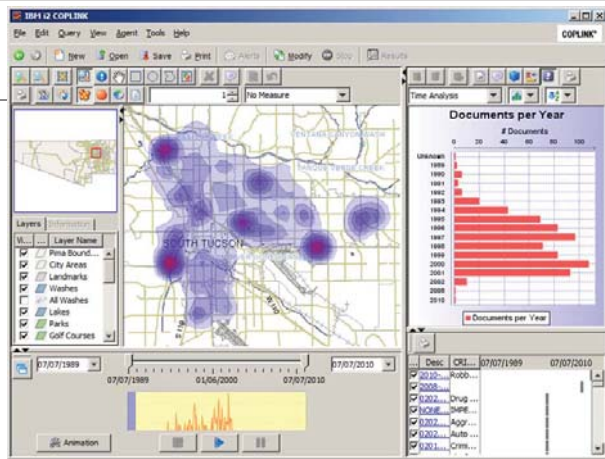
36



37



40



38

Social Media Analytics

- Social Media
 - Social interactions among people
 - People freely reveal their feelings
 - Highly dynamic

41

Example: US Law Enforcement

- Predictive policing
 - Predict time and location with high-probability of criminal activities
 - Prevent crimes before they happen
- Crime mapping
 - RAIDS Online
 - Connects law enforcement with the community to reduce crime and improve public safety

39

Different Types of Social Media

- Collaborative projects (e.g., Wikipedia, Wiktionary)
- Blogs and microblogs (e.g., Twitter, Weibo)
- Content communities (e.g., YouTube)
- Social networking sites (e.g., Facebook)
- Virtual game worlds (e.g., World of Warcraft)
- Virtual social worlds (e.g., Second Life)

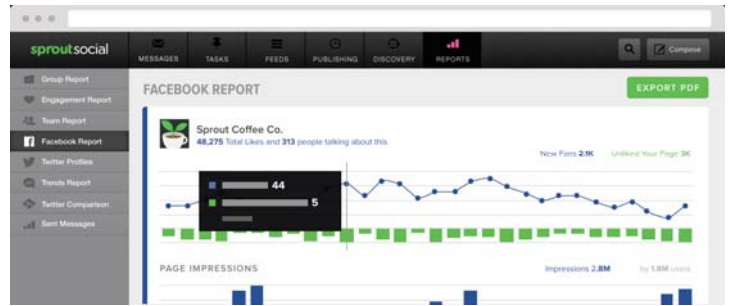
42

2015 Facebook usage

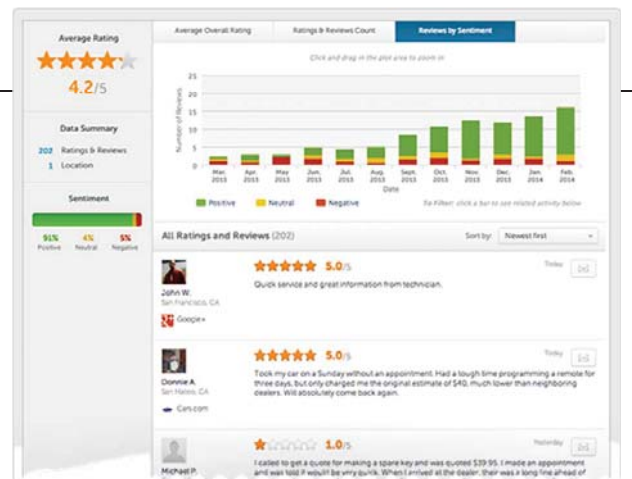
facebook Community Update

- 1.49 Billion** people on Facebook each month
- 1.5 Billion** searches daily
- 800 Million** people on WhatsApp each month
- 1 Billion** people offered access through Internet.org
- 700 Million** people on Messenger each month
- 850 Million** people using Groups on Facebook
- 450 Million** people using Events on Facebook
- 300 Million** people on Instagram each month
- 40 Million** small businesses using Pages

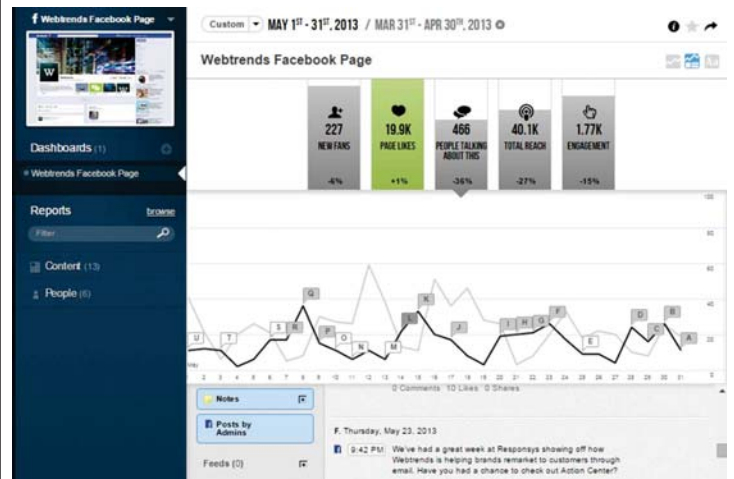
Social Media Analytics

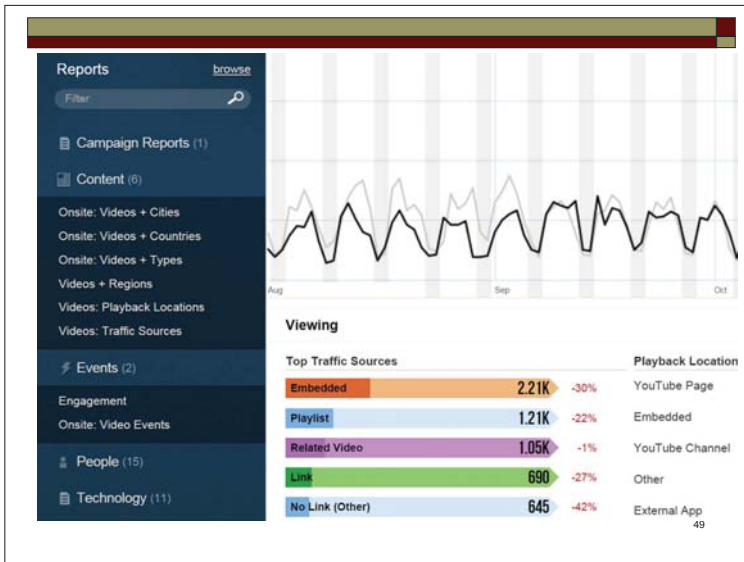


Social Media Analytics Tools and Vendors



Social Media Analytics





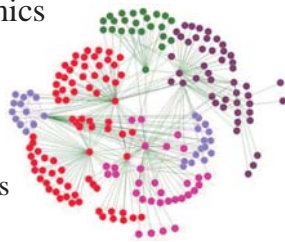
Application Case: HK Online Public Opinion

- Case Study: National education debate
- Controversy on the implementation of national education in secondary and primary schools
- Study period lasted from January 1, 2012 to May, 31, 2012.

52

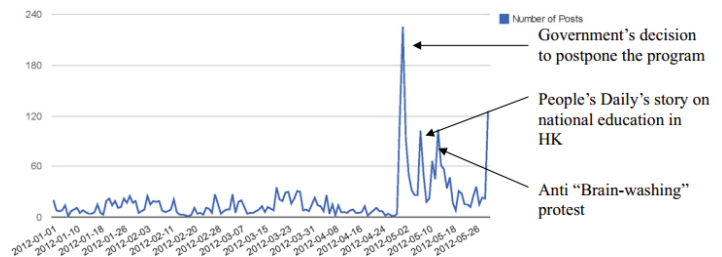
Social Network Analysis

- Social Network** - social structure composed of individuals linked to each other
- Analysis of social dynamics
- Identify
 - Opinion leaders
 - Bridges
 - Clusters and social circles



50

Time Trend of Total Number of Posts



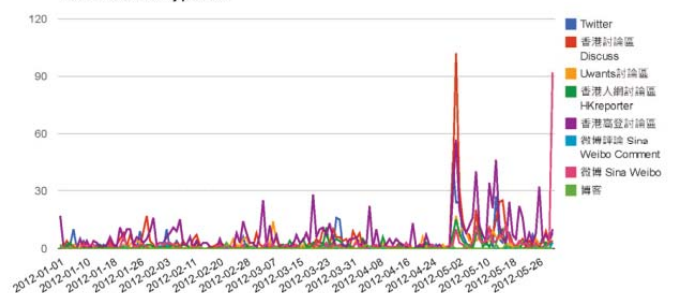
53

Application Case: HK Online Public Opinion

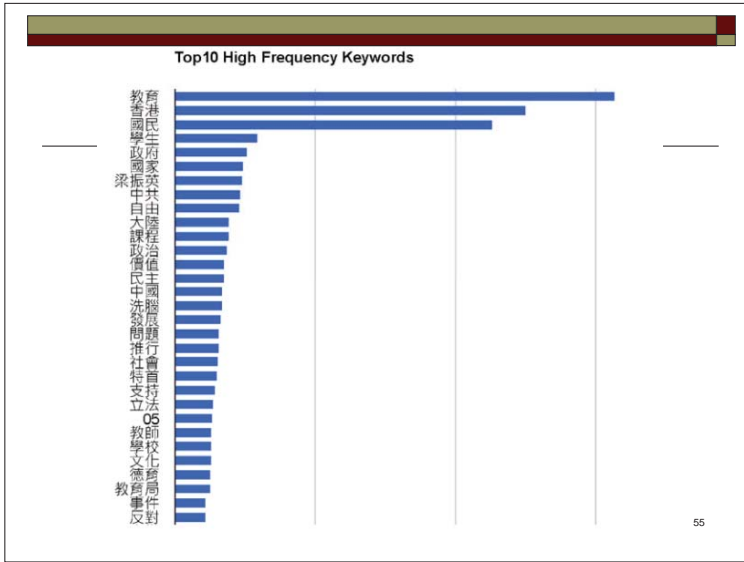
- Challenges
 - Online platforms become an important venue to understand public opinions
 - Too much information for analysis and monitoring
- Data collection from major Hong Kong based online opinion platforms
 - Discussion forums
 - Uwants, Discuss HK, Golden, HK Reporter
 - Twitter
 - Sina weibo
 - Blogs
 - Facebook pages, groups, and events

51

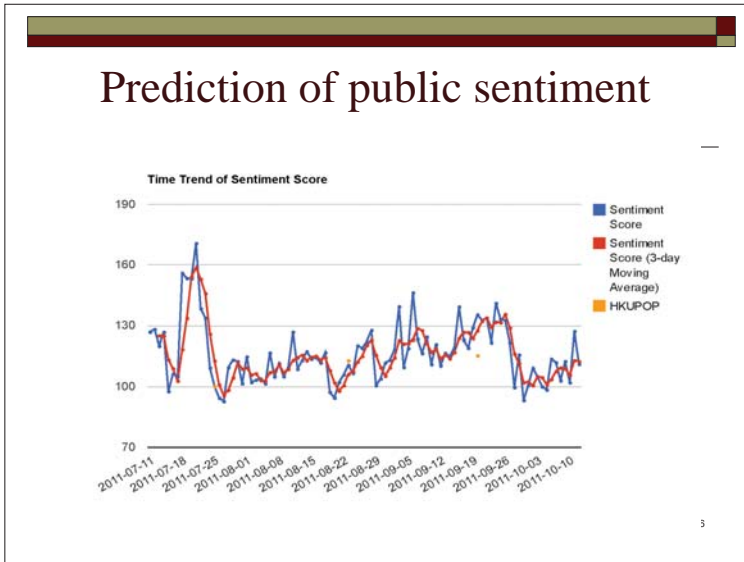
Time Trend of Posts by platforms



54



58



59

Application Case: Singapore Social Media Analytics

- Analyze and visualize social media
 - <http://research.larc.smu.edu.sg/palanteer/>
 - <http://research.larc.smu.edu.sg/palanteer/>

57

60

