# Introduction to Big Data

Based on "Big Data: Hype or Hallelujah?" by Elena Baralis
http://dbdmg.polito.it/wordpress/wp-content/uploads/2010/12/BigData_2015_2x.pdf

# Big data

# Google Flu trends

google.org Flu Trends

Google.org home
**Flu Trends**
Select country/region

Home
How does this work?
FAQ

**Flu activity**
Intense
High
Moderate
Low
Minimal

**Explore flu trends around the world**

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. Learn more »

Download world flu activity data

- February 2010
  - Google detected flu outbreak two weeks ahead of CDC data (Centers for Disease Control and Prevention – U.S.A)
  - Based on the analysis of Google search queries

8.866
6.650
4.433
2.217

2004  2005  2006  2007  2008  2009  2010  2011  2012  2013  2014

# Google Flu trends

google.org Flu Trends

Google.org home
**Flu Trends**
Select country/region

Home
How does this work?
FAQ

**Flu activity**
Intense
High
Moderate
Low
Minimal

**Explore flu trends around the world**
We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. Learn more »

- February 2010
  - Google detected flu outbreak two weeks ahead of CDC data (Centers for Disease Control and ...S.A)
  - ...alysis of ...ueries

# Nowcasting

8.866

6.650

4.433

2.217

2004  2005  2006  2007  2008  2009  2010  2011  2012  2013  2014

# Data on the Internet...

- Internet live stats
  - http://www.internetlivestats.com/

**4,485,508,861**
Internet Users in the world

**1,752,142,970**
Total number of Websites

**193,688,718,339**
Emails sent today

**5,222,289,027**
Google searches today

**4,990,992**
Blog posts written today

**572,159,945**
Tweets sent today

**5,348,093,035**
Videos viewed today
on YouTube

**62,832,046**
Photos uploaded today
on Instagram

**107,175,151**
Tumblr posts today

**2,435,900,914**
Facebook active users

**795,537,418**
Google+ active users

**357,398,865**
Twitter active users

**278,573,312**
Pinterest active users

**287,236,096**
Skype calls today

**110,341**
Websites hacked today

**5,664,059,486** GB
Internet traffic today

**3,065,544** MWh
Electricity used today
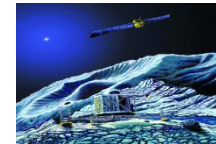for the Internet

**2,507,959** tons
$CO_2$ emissions today

5

# Who generates big data?

- User Generated Content (Web & Mobile)
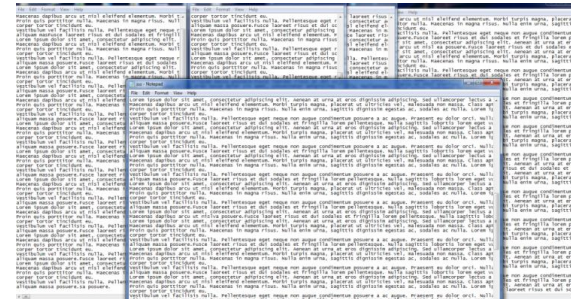  - E.g., Facebook, Instagram, Yelp, TripAdvisor, Twitter, YouTube

- Health and scientific computing

# Who generates big data?

- ## Log files
  - Web server log files, machine system log files

- ## Internet Of Things (IoT)
  - Sensor networks, RFIDs, smart meters

# An example of Big data at work

Crowdsourcing

Sensing

Map data

Computing

Real time traffic info
Travel time forecast/nowcast

# What is big data?



- Many different definitions
  - "Data whose scale, diversity and complexity require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"

# What is big data?



- Many different definitions
  - "Data whose **scale**, **diversity** and **complexity** require new architectures, techniques, algorithms and analytics to manage it and extract value and hidden knowledge from it"

# What is big data?



- Many different definitions
    - "Data whose scale, diversity and complexity require new **architectures**, **techniques**, **algorithms** and **analytics** to manage it and extract value and hidden knowledge from it"

# The Vs of big data

- The 3Vs of big data
  - **V**olume: scale of data
  - **V**ariety: different forms of data
  - **V**elocity: analysis of streaming data
- … but also
  - **V**eracity: uncertainty of data
  - **V**alue: exploit information provided by data

# The Vs of big data

- **V**olume

terabytes | petabytes | exabytes | zettabytes

*the amount of data stored by the average company today*

  - Data volume increases exponentially over time
  - 44x increase from 2009 to 2020
    - Digital data 35 ZB in 2020



Data storage growth



Twitter: Tweets Per Day



The Digital Universe 2009-2020

Growing By A Factor Of 44

2009: 0.8 Zb

2020: 35.2 Zettabytes

# The Vs of big data

- **V**ariety
  - Various formats, types and structures
    - Numerical data, image data, audio, video, text, time series



  - A single application may generate many different formats
    - Heterogeneous data
    - Complex data integration problem

# The Vs of big data

- **V**elocity
  - Fast data generation rate
    - Streaming data
  - Very fast data processing to ensure timeliness

# The Vs of big data

- **V**eracity
  - Data quality

# The Vs of big data

- **V**alue
  - Translate data into business advantage

Example: US economy

Size of bubble indicates relative contribution to GDP

High

Big data: ease-of-capture index[1]

Utilities
Natural resources
Manufacturing
Professional services
Accommodation and food
Construction
Administrative services
Other services

Health care providers
Computers and other electronic products
Information
Finance and insurance
Transportation and warehousing
Real estate
Management of companies
Wholesale trade
Retail trade

Educational services
Arts and entertainment
Government

Low

High

Big data: value potential index[1]

[1] For detailed explication of metrics, see appendix in McKinsey Global Institute full report *Big data: The next frontier for innovation, competition, and productivity*, available free of charge online at mckinsey.com/mgi.

Source: US Bureau of Labor Statistics; McKinsey Global Institute analysis

17

# Big data value chain

| Generation | Acquisition | Storage | Analysis |

- Generation
  - Passive recording
    - Typically structured data
    - Bank trading transactions, shopping records, government sector archives
  - Active generation
    - Semistructured or unstructured data
    - User-generated content, e.g., social networks
  - Automatic production
    - Location-aware, context-dependent, highly mobile data
    - Sensor-based Internet-enabled devices

# Big data value chain

Generation → **Acquisition** → Storage → Analysis

- Acquisition
  - Collection
    - Pull-based, e.g., web crawler
    - Push-based, e.g., video surveillance, click stream
  - Transmission
    - Transfer to data center over high capacity links
  - Preprocessing
    - Integration, cleaning, redundancy elimination

# Big data value chain

Generation ▸ Acquisition ▸ Storage ▸ Analysis

- Storage
  - Storage infrastructure
    - Storage technology, e.g., HDD, SSD
    - Networking architecture, e.g., DAS, NAS, SAN
  - Data management
    - File systems (HDFS), key-value stores (Memcached), column-oriented databases (Cassandra), document databases (MongoDB)
  - Programming models
    - Map reduce, stream processing, graph processing

# Big data value chain

Generation ▸ Acquisition ▸ Storage ▸ **Analysis**
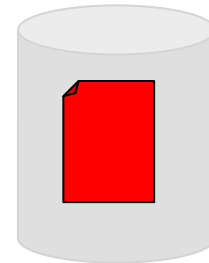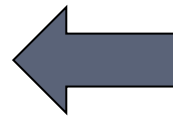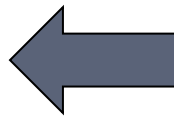
- Analysis
  - Objectives
    - Descriptive analytics, predictive analytics, prescriptive analytics
  - Methods
    - Statistical analysis, data mining, text mining, network and graph data mining
    - Clustering, classification and regression, association analysis
  - Diverse domains call for customized techniques

# Big data challenges

- Technology and infrastructure
  - New architectures, programming paradigms and techniques are needed
- Data management and analysis
  - New emphasis on "data"
  - ➡ **Data science**

# The bottleneck

- Processors process data
- Hard drives store data
- We need to transfer data from the disk to the processor

# The solution

- **Transfer the processing power to the data**
- Multiple distributed disks
  - Each one holding a portion of a large dataset
- Process in parallel different file portions from different disks